

Optimal Learning Rates for Localized SVMs

Mona Eberts

MONA.EBERTS@MATHEMATIK.UNI-STUTTGART.DE

Ingo Steinwart

INGO.STEINWART@MATHEMATIK.UNI-STUTTGART.DE

Institute for Stochastics and Applications

University of Stuttgart

70569 Stuttgart, Germany

Abstract

One of the limiting factors of using support vector machines (SVMs) in large scale applications are their super-linear computational requirements in terms of the number of training samples. To address this issue, several approaches that train SVMs on many small chunks of large data sets separately have been proposed in the literature. So far, however, almost all these approaches have only been empirically investigated. In addition, their motivation was always based on computational requirements. In this work, we consider a localized SVM approach based upon a partition of the input space. For this local SVM, we derive a general oracle inequality. Then we apply this oracle inequality to least squares regression using Gaussian kernels and deduce local learning rates that are essentially minimax optimal under some standard smoothness assumptions on the regression function. This gives the first motivation for using local SVMs that is not based on computational requirements but on theoretical predictions on the generalization performance. We further introduce a data-dependent parameter selection method for our local SVM approach and show that this method achieves the same learning rates as before. Finally, we present some larger scale experiments for our localized SVM showing that it achieves essentially the same test performance as a global SVM for a fraction of the computational requirements. In addition, it turns out that the computational requirements for the local SVMs are similar to those of a vanilla random chunk approach, while the achieved test errors are significantly better.

Keywords: least squares regression, support vector machines, localization

1. Introduction

Based on a training set $D := ((x_1, y_1), \dots, (x_n, y_n))$ of i.i.d. input/output observations drawn from an unknown distribution P on $X \times Y$, where $X \subset \mathbb{R}^d$ and $Y \subset \mathbb{R}$, the goal of non-parametric regression is to find a function $f_D : X \rightarrow \mathbb{R}$ such that important characteristics of the conditional distribution $P(Y|x)$, $x \in X$, can be recovered. For instance, an f_D approximating the conditional mean $\mathbb{E}(Y|x)$, $x \in X$, is sought in the non-parametric least squares regression. This classical non-parametric regression problem has been extensively studied in the literature, where a general reference is the book (Györfi et al., 2002), presenting plenty of results concerning the non-parametric least squares regression.

In the literature, there are many learning methods that solve the non-parametric regression problems, some of them are e.g. described in (Györfi et al., 2002; Koenker, 2005; Simonoff, 1996). In this paper, we utilize some kernel-based regularized empirical risk minimizers, also known as support vector machines (SVMs), which solve the regularized

problem

$$f_{D,\lambda} \in \arg \min_{f \in H} \lambda \|f\|_H^2 + \mathcal{R}_{L,D}(f) . \quad (1)$$

Here, $\lambda > 0$ is a fixed real number and H is a reproducing kernel Hilbert space (RKHS) over X with reproducing kernel $k : X \times X \rightarrow \mathbb{R}$, see e.g. (Aronszajn, 1950; Berlinet and Thomas-Agnan, 2004; Steinwart and Christmann, 2008a). Besides, $\mathcal{R}_{L,D}(f)$ denotes the empirical risk of a function $f : X \rightarrow \mathbb{R}$, that is

$$\mathcal{R}_{L,D}(f) = \frac{1}{n} \sum_{i=1}^n L(x_i, y_i, f(x_i)) ,$$

where D is the empirical measure associated to the data D defined by $D := \frac{1}{n} \sum_{i=1}^n \delta_{(x_i, y_i)}$ with Dirac measure $\delta_{(x_i, y_i)}$ at (x_i, y_i) . Note that the empirical SVM solution $f_{D,\lambda}$ exists and is unique (cf. Steinwart and Christmann, 2008a, Theorem 5.5) whenever the loss L is convex in its last argument. Moreover, an SVM is L -risk consistent under a few assumptions on the RKHS H and the regularization parameter λ , see (Steinwart and Christmann, 2008a, Section 6.4) for more details. Besides, it is worth mentioning that the ability to choose the RKHS H as well as the loss function L in (1) provides the possibility to flexibly apply SVMs to various learning problems. Namely, the learning target is modeled by the loss function, e.g. the least squares loss is used to estimate the conditional mean. Moreover, since RKHSs are defined on arbitrary X , data types that are not \mathbb{R}^d -valued can be handled, too. Furthermore, SVMs are enjoying great popularity, since they can be implemented and applied in a relatively simple way and only have a few free parameters that can usually be determined by cross validation.

An essential theoretical task, which has attracted many considerations, is the investigation of learning rates for SVMs. For example, such rates for SVMs using the least squares loss and generic kernels can be found in (Cucker and Smale, 2002; De Vito et al., 2005; Smale and Zhou, 2007; Caponnetto and De Vito, 2007; Mendelson and Neeman, 2010; Steinwart et al., 2009) and the references therein, while similar rates for SVMs using the pinball loss can be found in (Steinwart and Christmann, 2008b, 2011). At this point, we do not want to take a closer look at these results, instead we relegate to (Eberts and Steinwart, 2013), where a detailed discussion can be found. More important for our purposes is the fact that Eberts and Steinwart (2011, 2013) establish (essentially) asymptotically optimal learning rates for least squares SVMs (LS-SVMs) using Gaussian RBF kernels. More precisely, for a domain $X \subset B_{\ell_2^d}$, $Y := [-M, M]$ with $M > 0$, a distribution P on $X \times Y$ such that P_X has a bounded Lebesgue density on X , and for f^* contained in the Sobolev space $W_2^\alpha(P_X)$, $\alpha \in \mathbb{N}$, or in the Besov space $B_{2,\infty}^\alpha(P_X)$, $\alpha \geq 1$, respectively, the LS-SVM using Gaussian kernels learns for all $\xi > 0$ with rate $n^{-\frac{2\alpha}{2\alpha+d} + \xi}$ with a high probability. Although these rates are essentially asymptotically optimal, they depend on the order of smoothness of the regression function on the *entire* input space X . That is, if the regression function f^* is on some area of X smoother than on another area, the learning rate is determined by the part of X , where the regression function f^* is least smooth (cf. Figure 1).

In contrast to this, it would be desirable to achieve a learning rate on every region of X that corresponds with the order of smoothness of f^* on this region. Therefore, one of our goals of this paper is to modify the standard SVM approach such that we achieve local learning rates that are asymptotically optimal. Our technique to achieve such local learning rates is a special local SVM approach. Local SVMs have been extensively investigated in the literature to speed-up the training time, see for instance, the early works (Bottou and Vapnik, 1992; Vapnik and Bottou, 1993). The basic idea of many local approaches is to *a)* split the training data and just consider a few examples near a testing sample, *b)* train on this small subset of the training data, and *c)* use the solution for a prediction w.r.t. the

test sample. Here, many up-to-date investigations use SVMs to train on the local data set but, yet there are different ways to split the whole training data set into smaller, local sets. For example, Chang et al. (2010); Wu et al. (1999); Bennett and Blue (1998) use decision trees while in (Hable, 2013; Segata and Blanzieri, 2010, 2008; Blanzieri and Melgani, 2008; Blanzieri and Bryl, 2007a,b; Zhang et al., 2006) local subsets are built considering k nearest neighbors. The latter approaches further vary, for example, Zhang et al. (2006); Blanzieri and Bryl (2007a); Hable (2013) consider different metrics w.r.t. the input space whereas Segata and Blanzieri (2008); Blanzieri and Melgani (2008); Blanzieri and Bryl (2007b) consider metrics w.r.t. the feature space. Nonetheless, the basic idea of all these articles is that an SVM problem based on k training samples is solved for *each* test sample. Another approach using k nearest neighbors is investigated in (Segata and Blanzieri, 2010). Here, k -neighborhoods consisting of training samples and collectively covering the training data set are constructed and an SVM is calculated on each neighborhood. The prediction for a test sample is then made according to the nearest training sample that is a center of a k -neighborhood. As for the other nearest neighbor approaches, however, the results are mainly experimental. An exception to this rule is (Hable, 2013), where universal consistency for localized versions of SVMs, or more precisely, a large class of regularized kernel methods, is proven. Another article presenting theoretical results for localized versions of learning methods is (Zakai and Ritov, 2009). Here, the authors show that a consistent learning method behaves locally, i.e. the prediction is essentially influenced by close by samples. However, this result is based on a localization technique considering only training samples contained in a neighborhood with a fixed radius and center x when an estimate in x is sought. Probably closest to our approach is the one examined in (Cheng et al., 2010) and (Cheng et al., 2007), where the training data is splitted into clusters and then an SVM is trained on each cluster. However, the presented results are only of experimental character.

In this article, we partition the input space X according to a cover of X with radius r_n and build an SVM model for each partition cell. The following section is dedicated to the detailed description of this method. Section 3 then presents some theoretical results that

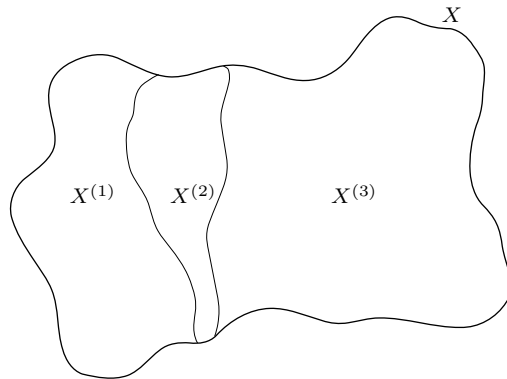


Figure 1: The input space X is partitioned by $X^{(1)}$, $X^{(2)}$, and $X^{(3)}$ such that the regression function f^* is less smooth on $X^{(2)}$ compared to $X^{(1)}$ and $X^{(3)}$. However, it is desirable to achieve locally optimal learning rates.

enable the analysis of this new method. For example, we examine extensions and direct sums of RKHSs. At the end of Section 3, we finally present a first oracle inequality for the localized SVM. In Section 4, we focus on RKHSs using Gaussian RBF kernels and, in conjunction with that, we study some entropy estimates. After that, Section 5 concentrates on the least squares loss and introduces an oracle inequality and learning rates for our localized SVM method using Gaussian kernels. Moreover, a data-dependent parameter selection method is studied that induces the same rates. Section 6 then presents some experimental results w.r.t. the localized SVM technique. All proofs can be found in Section 7, and the appendix contains various tables displaying detailed results of our experiments.

2. Description of the Localized SVM Approach

In this section, we introduce some general notations and assumptions. Based on the latter we modify the standard SVM approach. Let us start with the probability measure P on $X \times Y$, where $X \subset \mathbb{R}^d$ is non-empty and $Y := [-M, M]$ for some $M > 0$. Depending on the learning target one chooses a loss function L , i.e. a function $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ that is measurable. Then, for a measurable function $f : X \rightarrow \mathbb{R}$, the L -risk is defined by

$$\mathcal{R}_{L,P}(f) = \int_{X \times Y} L(x, y, f(x)) dP(x, y)$$

and the optimal L -risk, called the Bayes risk with respect to P and L , is given by

$$\mathcal{R}_{L,P}^* := \inf \{ \mathcal{R}_{L,P}(f) \mid f : X \rightarrow \mathbb{R} \text{ measurable} \}.$$

A measurable function $f_{L,P}^* : X \rightarrow \mathbb{R}$ with $\mathcal{R}_{L,P}(f_{L,P}^*) = \mathcal{R}_{L,P}^*$ is called a Bayes decision function. For the commonly used losses such as the least squares loss treated in Section 5 the Bayes decision function $f_{L,P}^*$ is P_X -almost surely $[-M, M]$ -valued, since $Y = [-M, M]$. In this case, it seems obvious to consider estimators with values in $[-M, M]$ on X . To this end, we now introduce the concept of clipping the decision function. Let \hat{t} be the clipped value of some $t \in \mathbb{R}$ at $\pm M$ defined by

$$\hat{t} := \begin{cases} -M & \text{if } t < -M \\ t & \text{if } t \in [-M, M] \\ M & \text{if } t > M. \end{cases}$$

Then a loss is called clippable at $M > 0$ if, for all $(x, y, t) \in X \times Y \times \mathbb{R}$, we have

$$L(x, y, \hat{t}) \leq L(x, y, t).$$

Obviously, the latter implies

$$\mathcal{R}_{L,P}(\hat{f}) \leq \mathcal{R}_{L,P}(f)$$

for all $f : X \rightarrow \mathbb{R}$. In other words, restricting the decision function to the interval $[-M, M]$ containing our labels cannot worsen the risk, in fact, clipping this function typically reduces the risk. Hence, we consider the clipped version \hat{f}_D of the decision function as well as the

risk $\mathcal{R}_{L,P}(\hat{f}_D)$ instead of the risk $\mathcal{R}_{L,P}(f_D)$ of the unclipped decision function. Note, this clipping idea does *not* change the learning method since it is performed *after* the training phase.

To modify the standard SVM approach (1), we assume that $(A_j)_{j=1,\dots,m}$ is a partition of X such that $A_j \neq \emptyset$ for every $j \in \{1, \dots, m\}$. Obviously, this implies $A_{j_1} \cap A_{j_2} = \emptyset$ for all $j_1, j_2 \in \{1, \dots, m\}$ with $j_1 \neq j_2$ and

$$X = \bigcup_{j=1}^m A_j.$$

Now, the basic idea of the approach developed in this paper is to consider for each set of the partition $(A_j)_{j=1,\dots,m}$ an individual SVM. To describe this approach in a mathematically rigorous way, we have to introduce some more definitions and notations. Let us begin with the index set

$$I_j := \{i \in \{1, \dots, n\} : x_i \in A_j\}, \quad j = 1, \dots, m,$$

indicating the samples of D contained in A_j , as well as the corresponding data set

$$D_j := \{(x_i, y_i) \in D : i \in I_j\}, \quad j = 1, \dots, m.$$

Moreover, for every $j \in \{1, \dots, m\}$, we define a (local) loss function $L_j : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ by

$$L_j(x, y, t) := \mathbb{1}_{A_j}(x) L(x, y, t), \quad (2)$$

where $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ is the loss that corresponds to our learning problem at hand. We further assume that H_j is an RKHS over A_j with kernel $k_j : A_j \times A_j \rightarrow \mathbb{R}$. Note that every function $f \in H_j$ is only defined on A_j even though a function $f_D : X \rightarrow \mathbb{R}$ is finally sought. To this end, for $f \in H_j$, we define a function $\hat{f} : X \rightarrow \mathbb{R}$ by

$$\hat{f}(x) := \begin{cases} f(x), & x \in A_j, \\ 0, & x \notin A_j. \end{cases}$$

Then the space $\hat{H}_j := \{\hat{f} : f \in H_j\}$ equipped with the norm

$$\|\hat{f}\|_{\hat{H}_j} := \|f\|_{H_j}, \quad \hat{f} \in \hat{H}_j,$$

is an RKHS on X (cf. Lemma 2). That is, \hat{H}_j is an isometrically isomorphic extension of the RKHS H_j on A_j to an RKHS on X . After all, we are now able to formulate a modified SVM approach. To this end, for every $j \in \{1, \dots, m\}$, consider the local SVM optimization problem

$$f_{D_j, \lambda_j} = \arg \min_{\hat{f} \in \hat{H}_j} \lambda_j \|\hat{f}\|_{\hat{H}_j}^2 + \frac{1}{n} \sum_{i=1}^n L_j(x_i, y_i, \hat{f}(x_i)), \quad (3)$$

where $\lambda_j > 0$ for every $j \in \{1, \dots, m\}$. Based on these empirical SVM solutions, we then define the decision function $f_{D,\lambda} : X \rightarrow \mathbb{R}$ by

$$f_{D,\lambda}(x) := \sum_{j=1}^m f_{D_j,\lambda_j}(x) = \sum_{j=1}^m \mathbb{1}_{A_j}(x) f_{D_j,\lambda_j}(x), \quad (4)$$

where $\lambda := (\lambda_1, \dots, \lambda_m)$. Here, clipping $f_{D,\lambda}$ at M yields

$$\hat{f}_{D,\lambda}(x) = \sum_{j=1}^m \mathbb{1}_{A_j}(x) \hat{f}_{D_j,\lambda_j}(x)$$

for every $x \in X$. Note that the empirical SVM solutions f_{D_j,λ_j} in (3) exist and are unique by (Steinwart and Christmann, 2008a, Theorem 5.5) and that, for arbitrary $j \in \{1, \dots, m\}$, $f_{D_j,\lambda_j} = 0$ if $x_i \notin A_j$ for all $i \in \{1, \dots, n\}$. In addition, the SVM optimization problem (3) equals the SVM optimization problem (1) using H_j , D_j , and the regularization parameter $\tilde{\lambda}_j := \frac{n}{|I_j|} \lambda_j$, since, for $\hat{f} \in \hat{H}_j$ and $f := \hat{f}|_{A_j}$, we have

$$\begin{aligned} \lambda_j \|\hat{f}\|_{\hat{H}_j}^2 + \frac{1}{n} \sum_{i=1}^n L_j(x_i, y_i, \hat{f}(x_i)) &= \lambda_j \|f\|_{H_j}^2 + \frac{1}{n} \sum_{i \in I_j} L(x_i, y_i, f(x_i)) \\ &= \frac{|I_j|}{n} \left(\tilde{\lambda}_j \|f\|_{H_j}^2 + \mathcal{R}_{L,D_j}(f) \right). \end{aligned}$$

That is, f_{D_j,λ_j} as in (3) and $h_{D_j,\tilde{\lambda}_j} := \arg \min_{f \in H_j} \tilde{\lambda}_j \|f\|_{H_j}^2 + \mathcal{R}_{L,D_j}(f)$ satisfy

$$h_{D_j,\tilde{\lambda}_j} = f_{D_j,\lambda_j}|_{A_j}.$$

For the sake of completeness, we briefly examine the Bayes risks w.r.t. P and L_j . To this end, let $X \subset \mathbb{R}^d$, $Y \subset \mathbb{R}$, $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a loss function and P be a distribution on $X \times Y$ such that a Bayes decision function $f_{L,P}^* : X \rightarrow \mathbb{R}$ exists. Then, for all $j \in \{1, \dots, m\}$ and losses L_j defined by (2), it is easy to show

$$\mathcal{R}_{L_j,P}(f_{L,P}^*) = \mathcal{R}_{L_j,P}^*,$$

whenever $f_{L,P}^*$ exists. In other words, a Bayes decision function $f_{L,P}^*$ w.r.t. P and L additionally is a Bayes decision function w.r.t. P and L_j . Moreover, for function spaces $\mathcal{F}_1, \dots, \mathcal{F}_m$ over X , we have

$$\sum_{j=1}^m \min_{f_j \in \mathcal{F}_j} \mathcal{R}_{L_j,D}(f_j) = \min_{f_1 \in \mathcal{F}_1, \dots, f_m \in \mathcal{F}_m} \sum_{j=1}^m \mathcal{R}_{L_j,D}(f_j) \quad (5)$$

by the construction of the loss L_j .

Let us now present an advantageous characteristic of our modified SVM, namely the required computing time. Solving an usual SVM problem has a computational cost of $\mathcal{O}(n^q)$ where $q \in [2, 3]$ and n is the sample size. For the new approach we consider m working sets of size n_1, \dots, n_m where $n_i \approx n_j$ for all $i, j \in \{1, \dots, m\}$, i.e. $n_i \approx \frac{n}{m}$. Then for each

working set an usual SVM problem has to be solved such that, altogether, the modified SVM induces a computational cost of $\mathcal{O}\left(m\left(\frac{n}{m}\right)^q\right)$. That is, for some $\beta > 0$ and $m \approx n^\beta$ our approach is computationally cheaper than a traditional SVM. Note that our strategy using a partition of the input space is a typical way to speed-up algorithms and handle large data sets. Other techniques that possess similar properties are e.g. applied in the articles cited in the introduction. Besides, we refer to (Tsang et al., 2007) and (Tsang et al., 2005) using enclosing ball problems to solve an SVM, to (Graf et al., 2005) presenting an model of multiple filtering SVMs and to (Collobert et al., 2001) investigating a mixture of SVMs based on several subsets of the training set.

To describe the above SVM approach $(A_j)_{j=1,\dots,m}$ only has to be some partition of X . However, for the theoretical investigations concerning learning rates of our new approach, we have to further specify the partition. To this end, we denote by $B_{\ell_2^d}$ the closed unit ball in the d -dimensional Euclidean space ℓ_2^d and we define balls B_1, \dots, B_m with radius $r > 0$ and mutually distinct centers $z_1, \dots, z_m \in X$ by

$$B_j := B_r(z_j) := \{x \in X : \|x - z_j\|_2 \leq r\}, \quad j \in \{1, \dots, m\}, \quad (6)$$

where $\|\cdot\|_2$ is the Euclidean norm in \mathbb{R}^d . Moreover, choose r and z_1, \dots, z_m such that

$$\bigcup_{j=1}^m B_j = X,$$

i.e. such that the balls B_1, \dots, B_m cover X (cf. Figure 2). The following well-known lemma relates the radius of such a cover with the number of centers.

Lemma 1 *Let $X \subset \mathbb{R}^d$ be a bounded subset, i.e. $X \subset cB_{\ell_2^d}$ for some constant $c > 0$. Then there exist balls $(B_j)_{j=1,\dots,m}$ with radius $r > 0$ covering X such that*

$$r \leq 8cm^{-\frac{1}{d}}.$$

For simplicity of notation, we assume in the following that $X \subset B_{\ell_2^d}$, i.e. according to Lemma 1 there exists a cover $(B_j)_{j=1,\dots,m}$ with

$$r \leq 8m^{-\frac{1}{d}}. \quad (7)$$

Finally, we can specify the partition $(A_j)_{j=1,\dots,m}$ of X by the following assumption.

- (A) Let $(A_j)_{j=1,\dots,m}$ be a partition of $X \subset B_{\ell_2^d}$ such that $\mathring{A}_j \neq \emptyset$ for every $j \in \{1, \dots, m\}$ and such that there exist mutually distinct $z_1, \dots, z_m \in X$ with $A_j \subset B_r(z_j) =: B_j$, where $(B_j)_{j=1,\dots,m}$ is a cover of X satisfying (7).

In the remaining sections we will frequently refer to Assumption (A). However, the results hold as well if we merely assume $z_1, \dots, z_m \in B_{\ell_2^d}$ instead of $z_1, \dots, z_m \in X \subset B_{\ell_2^d}$ in (A). The following example illustrates that (A) is indeed a natural assumption.

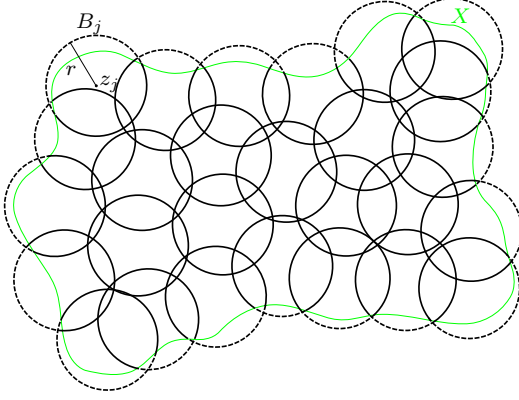


Figure 2: Cover $(B_j)_{j=1,\dots,m}$ of X , where B_1, \dots, B_m are balls with radius r and centers z_j ($j = 1, \dots, m$).

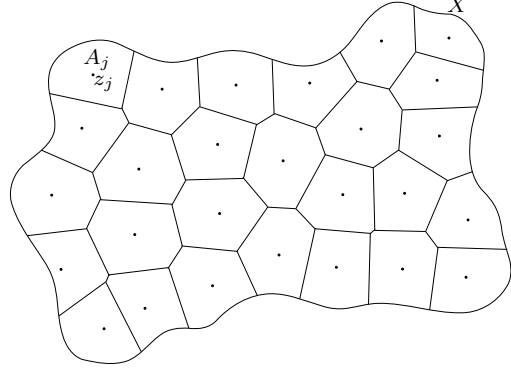


Figure 3: Voronoi partition $(A_j)_{j=1,\dots,m}$ of X defined by (8), where $A_j \subset B_j$ for every $j \in \{1, \dots, m\}$.

Example 1 For some $r > 0$, let us consider an r -net z_1, \dots, z_m of X , where z_1, \dots, z_m are mutually distinct. Based on these z_1, \dots, z_m , a Voronoi partition $(A_j)_{j=1,\dots,m}$ of X is defined by

$$A_j := \left\{ x \in X : j = \min \left\{ \arg \min_{k \in \{1, \dots, m\}} \|x - z_k\|_2 \right\} \right\}, \quad (8)$$

cf. Figure 3. That is, A_j contains all $x \in X$ such that the center z_j is the nearest center to x , and if there exist j_1 and j_2 with $j_1 < j_2$ and

$$\|x - z_{j_1}\|_2 = \|x - z_{j_2}\|_2 < \|x - z_k\|_2$$

for all $k \in \{1, \dots, m\} \setminus \{j_1, j_2\}$, then $x \in A_{j_1}$ since $j_1 < j_2$. In other words, they are resolved in favor of the smallest index of the involved centers. Moreover, it is obvious that $A_j \neq \emptyset$, $A_j \subset B_r(z_j)$ for all $j \in \{1, \dots, m\}$, $A_{j_1} \cap A_{j_2} = \emptyset$ for all $j_1, j_2 \in \{1, \dots, m\}$ with $j_1 \neq j_2$, and $X = \bigcup_{j=1}^m A_j$. In other words, a Voronoi partition based on an r -net z_1, \dots, z_m of X satisfies condition (A), if r and m fulfill (7).

Following Example 1, we call the learning method producing $f_{D,\lambda}$ given by (4) a *Voronoi partition support vector machine*, in short VP-SVM. Nevertheless, we just take a partition $(A_j)_{j=1,\dots,m}$ satisfying (A) as basis here instead of requesting $(A_j)_{j=1,\dots,m}$ to be a Voronoi partition.

Recall that our goal is to derive not only global but also local learning rates for this VP-SVM approach. To this end, we additionally consider an arbitrary measurable set $T \subset X$ such that $P_X(T) > 0$. Then we examine the learning rate of the VP-SVM on this subset T of X . To formalize this, it is necessary to introduce some basic notations related to T . Let us define the index set J_T by

$$J_T := \{j \in \{1, \dots, m\} : A_j \cap T \neq \emptyset\} \quad (9)$$

specifying every set A_j that has at least one common point with T . Note that, for every non-empty set $T \subset X$, the index set J_T is non-empty, too, i.e. $|J_T| \geq 1$. Besides, deriving

local rates on T requires us to investigate the excess risk of the VP-SVM with respect to the distribution P and the loss $L_T : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ defined by

$$L_T(x, y, t) := \mathbb{1}_T(x) L(x, y, t). \quad (10)$$

However, to manage the analysis we additionally need the loss $L_{J_T} : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ given by

$$L_{J_T}(x, y, t) := \mathbb{1}_{\bigcup_{j \in J_T} A_j}(x) L(x, y, t) \quad (11)$$

which may only be nonzero, if x is contained in some set A_j with $j \in J_T$. Note that the risks $\mathcal{R}_{L_T, P}(f)$ and $\mathcal{R}_{L_{J_T}, P}(f)$ quantify the quality of some function f just on T and

$$A_T := \bigcup_{j \in J_T} A_j \supset T,$$

respectively. Hence, examining the excess risks

$$\mathcal{R}_{L_T, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L_T, P}^* \leq \mathcal{R}_{L_{J_T}, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L_{J_T}, P}^*$$

leads to learning rates on A_T and implicitly on

T . Recapitulatory, let us declare a second set of assumptions.

(T) For $T \subset X$, we define an index set J_T by (9), loss functions $L_T, L_{J_T} : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ by (10) and (11), and the set $A_T := \bigcup_{j \in J_T} A_j$.

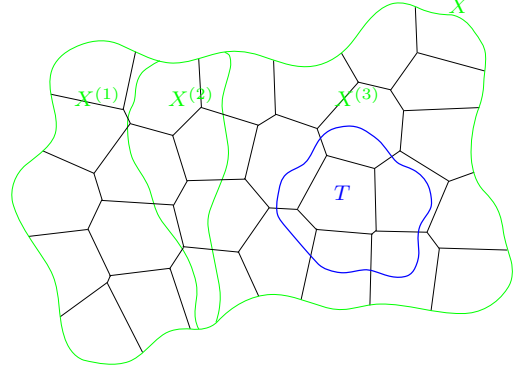


Figure 4: The input space X with the corresponding partition $(A_j)_{j=1, \dots, m}$ and the subset T , where the local learning rate should be examined.

3. An Oracle Inequality for VP-SVMs

In this section, we first focus on RKHSs and direct sums of RKHSs. Then we present a lemma that relates the risk of a function w.r.t. the general loss L to the risks w.r.t. the losses L_j . Finally, we establish a first oracle inequality for VP-SVMs.

Let us begin with some basic notations. For $q \in [1, \infty]$ and a measure ν , we denote by $L_q(\nu)$ the Lebesgue spaces of order q w.r.t. ν and for the Lebesgue measure μ on $X \subset \mathbb{R}^d$ we write $L_q(X) := L_q(\mu)$. In addition, for a measurable space X , the set of all real-valued measurable functions on X is given by $\mathcal{L}_0(X) := \{f : X \rightarrow \mathbb{R} \mid f \text{ measurable}\}$. Moreover, for a measure ν on X and measurable $\tilde{X} \subset X$, we define the trace measure $\nu|_{\tilde{X}}$ of ν in \tilde{X} by $\nu|_{\tilde{X}}(A) = \nu(A \cap \tilde{X})$ for every $A \subset X$.

Our first goal is to show that $\hat{f}_{D, \lambda}$ in (4) is actually an ordinary SVM solution. To this end, we consider an RKHS on some $A \subsetneq X$ and extend it to an RKHS on X by the following lemma, where we omit the obvious proof.

Lemma 2 *Let $A \subset X$ and H_A be an RKHS on A with corresponding kernel k_A . Denote by \hat{f} the extension of $f \in H_A$ to X defined by*

$$\hat{f}(x) := \begin{cases} f(x), & \text{for } x \in A, \\ 0, & \text{for } x \in X \setminus A. \end{cases}$$

Then the space $\hat{H}_A := \{\hat{f} : f \in H_A\}$ equipped with the norm

$$\|\hat{f}\|_{\hat{H}_A} := \|f\|_{H_A}$$

is an RKHS on X and its reproducing kernel is given by

$$\hat{k}_A(x, x') := \begin{cases} k_A(x, x'), & \text{if } x, x' \in A, \\ 0, & \text{else.} \end{cases} \quad (12)$$

Based on this lemma, we are now able to construct an RKHS by a direct sum of RKHSs \hat{H}_A and \hat{H}_B with $A, B \subset X$ and $A \cap B = \emptyset$. Here, we skip the proof once more, since the assertion follows immediately using, for example, orthonormal bases of \hat{H}_A and \hat{H}_B .

Lemma 3 *For $A, B \subset X$ such that $A \cap B = \emptyset$ and $A \cup B \subset X$, let H_A and H_B be RKHSs of k_A and k_B over A and B , respectively. Furthermore, let \hat{H}_A and \hat{H}_B be the RKHSs of all functions of H_A and H_B extended to X in the sense of Lemma 2 and let \hat{k}_A and \hat{k}_B given by (12) be the associated reproducing kernels. Then $\hat{H}_A \cap \hat{H}_B = \{0\}$ and hence the direct sum*

$$H := \hat{H}_A \oplus \hat{H}_B \quad (13)$$

exists. For $\lambda_A, \lambda_B > 0$ and $f \in H$, let $\hat{f}_A \in \hat{H}_A$ and $\hat{f}_B \in \hat{H}_B$ be the unique functions such that $f = \hat{f}_A + \hat{f}_B$. Then we define the norm $\|\cdot\|_H$ by

$$\|f\|_H^2 := \lambda_A \|\hat{f}_A\|_{\hat{H}_A}^2 + \lambda_B \|\hat{f}_B\|_{\hat{H}_B}^2 \quad (14)$$

and H equipped with the norm $\|\cdot\|_H$ is again an RKHS for which

$$k(x, x') := \lambda_A^{-1} \hat{k}_A(x, x') + \lambda_B^{-1} \hat{k}_B(x, x'), \quad x, x' \in X,$$

is the reproducing kernel.

To relate Lemmas 2 and 3 with (4), we have to introduce some more notations. For pairwise disjoint sets $A_1, \dots, A_m \subset X$, let H_j be an RKHS on A_j for every $j \in \{1, \dots, m\}$. Then, based on RKHSs $\hat{H}_1, \dots, \hat{H}_m$ on X defined by Lemma 2, the joined RKHSs can be designed analogously to Lemma 3. That is, for an arbitrary index set $J \subset \{1, \dots, m\}$ and a vector $\lambda = (\lambda_j)_{j \in J} \in (0, \infty)^{|J|}$, the direct sum

$$H_J := \bigoplus_{j \in J} \hat{H}_j = \left\{ f = \sum_{j \in J} f_j : f_j \in \hat{H}_j \text{ for all } j \in J \right\}$$

is again an RKHS equipped with the norm

$$\|f\|_{H_J}^2 = \sum_{j \in J} \lambda_j \|f_j\|_{\hat{H}_j}^2. \quad (15)$$

If $J = \{1, \dots, m\}$ we simply write

$$H := H_J \quad (16)$$

Note that H contains inter alia $f_{D, \lambda}$ given by (4). Summarizing, we can define another assumption set.

(H) For every $j \in \{1, \dots, m\}$, let H_j be a separable RKHS of measurable kernels k_j over A_j , where $A_1, \dots, A_m \subset X$ are pairwise disjoint and

$$\|k_j\|_{L_2(P_{X|A_j})}^2 := \int_X k_j(x, x) dP_{X|A_j}(x) < \infty.$$

Then we define RKHSs $\hat{H}_1, \dots, \hat{H}_m$ by Lemma 2 and the joined RKHS H by (16) equipped with the norm (15) for fixed $\lambda_1, \dots, \lambda_m > 0$.

Having designed a joined RKHS as above, a crucial property of its function's risks is expressed by the following lemma.

Lemma 4 *Let P be a distribution on $X \times Y$ and $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a loss function. For $A, B \subset X$ such that $A \cup B = X$ and $A \cap B = \emptyset$, define loss functions $L_A, L_B : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ by $L_A(x, y, t) = \mathbb{1}_A(x)L(x, y, t)$ and $L_B(x, y, t) = \mathbb{1}_B(x)L(x, y, t)$, respectively. Furthermore, let $f_A : X \rightarrow \mathbb{R}$ as well as $f_B : X \rightarrow \mathbb{R}$ be measurable functions and $f : X \rightarrow \mathbb{R}$ be defined by $f(x) = \mathbb{1}_A(x)f_A(x) + \mathbb{1}_B(x)f_B(x)$ for all $x \in X$. Then we have*

$$\mathcal{R}_{L,P}(f) = \mathcal{R}_{L_A,P}(f_A) + \mathcal{R}_{L_B,P}(f_B).$$

as well as

$$\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* = (\mathcal{R}_{L_A,P}(f_A) - \mathcal{R}_{L_A,P}^*) + (\mathcal{R}_{L_B,P}(f_B) - \mathcal{R}_{L_B,P}^*).$$

Note that Lemma 4 can be transferred to finite, pairwise disjoint unions. To be more precise, let us consider an arbitrary index set $J \subset \{1, \dots, m\}$ and define the corresponding loss function $L_J : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ by

$$L_J(x, y, t) := \mathbb{1}_{\bigcup_{j \in J} A_j}(x)L(x, y, t).$$

Now, it is straightforward to show

$$\mathcal{R}_{L_J,P}(f) = \sum_{j \in J} \mathcal{R}_{L_j,P}(f)$$

for every function $f : X \rightarrow \mathbb{R}$. Based on this generalization and the whole index set $J = \{1, \dots, m\}$, let us briefly consider Lemma 4 for the empirical measure D and for $f_{D,\lambda} = \sum_{j=1}^n \mathbb{1}_{A_j} f_{D_j,\lambda_j}$, where f_{D_j,λ_j} , $j = 1, \dots, m$, are defined by (3). Then, for an arbitrary $f \in H$, it immediately follows

$$\begin{aligned} \|f_{D,\lambda}\|_H^2 + \mathcal{R}_{L,D}(\hat{f}_{D,\lambda}) &= \sum_{j=1}^m \left(\lambda_j \|f_{D_j,\lambda_j}\|_{\hat{H}_j}^2 + \mathcal{R}_{L_j,D}(\hat{f}_{D,\lambda}) \right) \\ &\leq \sum_{j=1}^m \left(\lambda_j \|\mathbb{1}_{A_j} f\|_{\hat{H}_j}^2 + \mathcal{R}_{L_j,D}(f) \right) \\ &= \|f\|_H^2 + \mathcal{R}_{L,D}(f). \end{aligned} \tag{17}$$

That is, $f_{D,\lambda}$ is the decision function of an SVM using H and L as well as the regularization parameter $\lambda = 1$. In other words, the latter SVM equals the VP-SVM given by (4). This will be a key insight used in our analysis.

To derive an oracle inequality, i.e. an appropriate upper bound for the excess risk $\mathcal{R}_{L,J,P}(\widehat{f}_{D,\lambda}) - \mathcal{R}_{L,J,P}^*$ for some index set $J \subset \{1, \dots, m\}$, we have to introduce a few more notations. Let P be a distribution on $X \times Y$ such that a Bayes decision function $f_{L,P}^* : X \rightarrow [-M, M]$ exists, for some constant $M > 0$ at which L can be clipped. Moreover, we denote by $L \circ f$ the function $(x, y) \mapsto L(x, y, f(x))$. If there exist constants $B > 0$, $\vartheta \in [0, 1]$, and $V \geq B^{2-\vartheta}$ such that we have

$$L(x, y, t) \leq B, \quad (18)$$

$$\mathbb{E}_P (L \circ f - L \circ f_{L,P}^*)^2 \leq V \cdot (\mathbb{E}_P (L \circ f - L \circ f_{L,P}^*))^\vartheta, \quad (19)$$

for all $(x, y) \in X \times Y$, $t \in [-M, M]$, and $f : X \rightarrow [-M, M]$, we say that the supremum bound (18) and the variance bound (19), respectively, is fulfilled. Actually, (18) immediately yields

$$L_J(x, y, t) = \mathbb{1}_{\bigcup_{j \in J} A_j}(x) L(x, y, t) \leq L(x, y, t) \leq B$$

for all $(x, y) \in X \times Y$ and $t \in [-M, M]$, i.e. the supremum bound is also satisfied for L_J . Moreover, if (19) holds for all $f : X \rightarrow [-M, M]$, the variance bound using the loss L_J is satisfied, too. Indeed, by the use of $\tilde{f}(x) := \mathbb{1}_{\bigcup_{j \in J} A_j}(x) f(x) + \mathbb{1}_{X \setminus (\bigcup_{j \in J} A_j)}(x) f_{L,P}^*(x)$ for all $x \in X$, we have

$$\begin{aligned} \mathbb{E}_P (L_J \circ f - L_J \circ f_{L,P}^*)^2 &= \mathbb{E}_P (L_J \circ \tilde{f} - L_J \circ f_{L,P}^*)^2 \\ &= \mathbb{E}_P (L \circ \tilde{f} - L \circ f_{L,P}^*)^2 \\ &\leq V \cdot (\mathbb{E}_P (L \circ \tilde{f} - L \circ f_{L,P}^*))^\vartheta \\ &\leq V \cdot (\mathbb{E}_P (L_J \circ f - L_J \circ f_{L,P}^*))^\vartheta \end{aligned}$$

for all $f : X \rightarrow [-M, M]$. Let us quickly define a third assumption set.

(P) Let P be a distribution on $X \times Y$ such that the variance bound (19) is satisfied for constants $\vartheta \in [0, 1]$, $V \geq B^{2-\vartheta}$, and all functions $f : X \rightarrow [-M, M]$.

Up to now, there is still missing a classical tool that is used to derive learning rates, namely entropy numbers, see (Carl and Stephani, 1990) or (Steinwart and Christmann, 2008a, Definition A.5.26). Recall that, for normed spaces $(E, \|\cdot\|_E)$ and $(F, \|\cdot\|_F)$ as well as an integer $i \geq 1$, the i -th (dyadic) entropy number of a bounded, linear operator $S : E \rightarrow F$ is defined by

$$\begin{aligned} e_i(S : E \rightarrow F) &:= e_i(SB_E, \|\cdot\|_F) \\ &:= \inf \left\{ \varepsilon > 0 : \exists s_1, \dots, s_{2^{i-1}} \in SB_E \text{ such that } SB_E \subset \bigcup_{j=1}^{2^{i-1}} (s_j + \varepsilon B_F) \right\}, \end{aligned}$$

where we use the convention $\inf \emptyset := \infty$, and B_E as well as B_F denote the closed unit balls in E and F , respectively. Finally, we present a first oracle inequality involving an upper bound for the excess risk $\mathcal{R}_{L_J, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L_J, P}^*$, where $J \subset \{1, \dots, m\}$ is an arbitrary index set.

Theorem 5 *Let $L : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ be a locally Lipschitz continuous loss that can be clipped at $M > 0$ and that satisfies the supremum bound (18) for some $B > 0$. Based on a partition $(A_j)_{j=1, \dots, m}$ of X , where $A_j \neq \emptyset$ for every $j \in \{1, \dots, m\}$, we assume (H). Furthermore, for an arbitrary index set $J \subset \{1, \dots, m\}$, we suppose (P). Assume that, for fixed $n \geq 1$, there exist constants $p \in (0, 1)$ and $a_1, \dots, a_m > 0$ such that for all $j \in \{1, \dots, m\}$*

$$e_i(\text{id} : H_j \rightarrow L_2(P_{X|A_j})) \leq a_j i^{-\frac{1}{2p}}, \quad i \geq 1. \quad (20)$$

Finally, fix an $f_0 \in H$ and a constant $B_0 \geq B$ such that $\|L_J \circ f_0\|_\infty \leq B_0$. Then, for all fixed $\tau > 0$, $\lambda = (\lambda_1, \dots, \lambda_m) > 0$, and

$$a := \max \left\{ c_p \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}}, B \right\},$$

the VP-SVM given by (4) using $\hat{H}_1, \dots, \hat{H}_m$ and L_J satisfies

$$\begin{aligned} & \sum_{j=1}^m \lambda_j \|f_{D_j, \lambda_j}\|_{\hat{H}_j}^2 + \mathcal{R}_{L_J, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L_J, P}^* \\ & \leq 9 \left(\sum_{j=1}^m \lambda_j \|\mathbb{1}_{A_j} f_0\|_{\hat{H}_j}^2 + \mathcal{R}_{L_J, P}(f_0) - \mathcal{R}_{L_J, P}^* \right) + C(a^{2p} n^{-1})^{\frac{1}{2-p-\vartheta+p}} + 3 \left(\frac{72V\tau}{n} \right)^{\frac{1}{2-\vartheta}} + \frac{15B_0\tau}{n} \end{aligned}$$

with probability P^n not less than $1 - 3e^{-\tau}$, where $C > 0$ is a constant only depending on p, M, V, ϑ , and B .

The above theorem deals with the case of a partition with quite a few sets A_j , $j \in \{1, \dots, m\}$. However, if we consider a partition consisting of just one set A_1 , i.e. $A_1 = X$, Theorem 5 is supposed to provide an oracle inequality that is comparable to the already known ones. To make that sure, let us briefly consider the case $m = 1$ and hence $A_1 := X$, $\lambda_1 := \lambda$ as well as RKHSs $H_1 = \hat{H}_1 = H$ over X with $\|\cdot\|_H^2 = \lambda \|\cdot\|_{H_1}^2$. Note that in this case we have $f_{D, \lambda} = f_{D_1, \lambda_1}$. If (20) holds for H_1 , Theorem 5 yields that an SVM using H and $L_J = L$ satisfies

$$\begin{aligned} & \lambda \|f_{D, \lambda}\|_{H_1}^2 + \mathcal{R}_{L, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L, P}^* \\ & \leq 9 (\lambda \|f_0\|_{H_1}^2 + \mathcal{R}_{L, P}(f_0) - \mathcal{R}_{L, P}^*) + C \left(\tilde{c}_p \frac{a_1^{2p}}{\lambda^p n} \right)^{\frac{1}{2-p-\vartheta+p}} + 3 \left(\frac{72V\tau}{n} \right)^{\frac{1}{2-\vartheta}} + \frac{15B_0\tau}{n} \end{aligned}$$

with probability P^n not less than $1 - 3e^{-\tau}$ for fixed $\tau > 0$. Note that this oracle inequality indeed matches with the one stated in (Steinwart and Christmann, 2008a, Theorem 7.23) apart from the constant \tilde{c}_p , which is, however, only depending on p , ϑ , and B .

In the following section, we focus on RKHSs using Gaussian RBF kernels and examine the associated entropy numbers to specify (20). Subsequently in Section 5, we additionally consider the least squares loss and adapt the oracle inequality of Theorem 5.

4. Entropy Estimates for Local Gaussian RKHSs

In this section, we refine assumption (20). More precisely, in the subsequent theorem we determine an upper bound for the entropy numbers of the operator $\text{id} : H_\gamma(A) \rightarrow L_2(P_{X|A})$, where $H_\gamma(A)$ is the RKHS over A of the Gaussian RBF kernel k_γ on $A \subset \mathbb{R}^d$ defined by

$$k_\gamma(x, x') := \exp(-\gamma^{-2} \|x - x'\|_2^2), \quad x, x' \in A,$$

for some width $\gamma > 0$.

Theorem 6 *Let $X \subset \mathbb{R}^d$, P_X be a distribution on X and $A \subset X$ be such that $\mathring{A} \neq \emptyset$ and such that there exists an Euclidean ball $B \subset \mathbb{R}^d$ with radius $r > 0$ containing A , i.e. $A \subset B$. Moreover, for $0 < \gamma \leq r$, let $H_\gamma(A)$ be the RKHS of the Gaussian RBF kernel k_γ over A . Then, for all $p \in (0, 1)$, there exists a constant $c_p > 0$ such that*

$$e_i(\text{id} : H_\gamma(A) \rightarrow L_2(P_{X|A})) \leq c_p \sqrt{P_X(A)} r^{\frac{d+2p}{2p}} \gamma^{-\frac{d+2p}{2p}} i^{-\frac{1}{2p}}, \quad i \geq 1.$$

Obviously, this theorem specifies assumption (20). Now, for the Gaussian case we elaborate assumption (H) and introduce the following additional set of assumptions.

(G) Let A_1, \dots, A_m be pairwise disjoint subsets of X with non-empty interior such that, for some fixed $r > 0$ and every $j \in \{1, \dots, m\}$, $\sup_{x, x' \in A_j} \|x - x'\|_2 \leq 2r$ is satisfied. Furthermore, for every $j \in \{1, \dots, m\}$, let $H_j := H_{\gamma_j}(A_j)$ be the RKHS of the Gaussian kernel k_{γ_j} with width $\gamma_j \in (0, r]$ over A_j . Consequently, for $\lambda := (\lambda_1, \dots, \lambda_m) \in (0, \infty)^m$, we define the joined RKHS $H := \bigoplus_{j=1}^m \hat{H}_{\gamma_j}(A_j)$ by (16) equipped with the norm (15).

Since we do not consider SVMs with a fixed kernel, we use a more detailed notation than (3) and (4) in the following specifying the kernel width γ_j of the RKHS $H_{\gamma_j}(A_j)$ at hand. For all $j \in \{1, \dots, m\}$ and $\gamma := (\gamma_1, \dots, \gamma_m)$, we thus write

$$f_{D_j, \lambda_j, \gamma_j} = \arg \min_{f \in \hat{H}_{\gamma_j}(A_j)} \lambda_j \|f\|_{\hat{H}_{\gamma_j}(A_j)}^2 + \frac{1}{n} \sum_{i=1}^n L_j(x_i, y_i, f(x_i)),$$

and

$$f_{D, \lambda, \gamma} := \sum_{j=1}^m f_{D_j, \lambda_j, \gamma_j}$$

instead of f_{D_j, λ_j} and $f_{D, \lambda}$ in the remainder of this work.

In the subsequent section, we consider the least squares loss which, together with Assumption (G) and Theorem 6, allows us to elaborate the oracle inequality stated in Theorem 5 so that we finally obtain learning rates.

5. Learning Rates for Least Squares VP-SVMs

In this section, the non-parametric least squares regression problem is considered using the least squares loss $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ defined by $L(y, t) := (y - t)^2$. It is well known that, in this case, the Bayes decision function $f_{L,P}^* : \mathbb{R}^d \rightarrow \mathbb{R}$ is given by $f_{L,P}^*(x) = \mathbb{E}_P(Y|x)$ for P_X -almost all $x \in \mathbb{R}^d$. Moreover, this function is unique up to zero-sets. Besides, for the least squares loss the equality

$$\mathcal{R}_{L,P}(f) - \mathcal{R}_{L,P}^* = \|f - f_{L,P}^*\|_{L_2(P_X)}^2$$

can be shown by some simple, well-known transformations. Recall that T is a non-empty subset of X , where the index set J_T defined by (9) indicates every set A_j of the partition $(A_j)_{j=1,\dots,m}$ of X that shares at least one point with T . The associated loss function $L_{J_T} : X \times Y \times \mathbb{R} \rightarrow [0, \infty)$ is defined by (11).

5.1 Basic Oracle Inequalities for LS-VP-SVMs

To formulate oracle inequalities and derive rates for VP-SVMs using the least squares loss, the target function $f_{L,P}^*$ is assumed to satisfy certain smoothness conditions. To this end, we initially recall the modulus of smoothness, a device to measure the smoothness of functions (see e.g. (Devore and Lorentz, 1993, p. 44), (Devore and Popov, 1988, p. 398), and (Berens and DeVore, 1978, p. 360)). Denote by $\|\cdot\|_2$ the Euclidean norm and let $X \subset \mathbb{R}^d$ be a subset with non-empty interior, ν be an arbitrary measure on X , $p \in (0, \infty]$, and $f : X \rightarrow \mathbb{R}$ be contained in $L_p(\nu)$. Then, for $s \in \mathbb{N}$, the s -th modulus of smoothness of f is defined by

$$\omega_{s,L_p(\nu)}(f, t) = \sup_{\|h\|_2 \leq t} \|\Delta_h^s(f, \cdot)\|_{L_p(\nu)} \quad , \quad t \geq 0 \quad ,$$

where $\Delta_h^s(f, \cdot)$ denotes the s -th difference of f given by

$$\Delta_h^s(f, x) = \begin{cases} \sum_{j=0}^s \binom{s}{j} (-1)^{s-j} f(x + jh) & \text{if } x \in X_{s,h} \\ 0 & \text{if } x \notin X_{s,h} \end{cases}$$

for $h = (h_1, \dots, h_d) \in [0, \infty)^d$ and $X_{s,h} := \{x \in X : x + th \in X \text{ f.a. } t \in [0, s]\}$. Based on the modulus of smoothness, we introduce Besov spaces, i.e. function spaces that provide a finer scale of smoothness than the commonly used Sobolev spaces and that will thus be assumed to contain the target function later on. To this end, let $1 \leq p, q \leq \infty$, $\alpha > 0$, $s := \lfloor \alpha \rfloor + 1$, and ν be an arbitrary measure. Then the Besov space $B_{p,q}^\alpha(\nu)$ is defined by

$$B_{p,q}^\alpha(\nu) := \left\{ f \in L_p(\nu) : |f|_{B_{p,q}^\alpha(\nu)} < \infty \right\} \quad ,$$

where the seminorm $| \cdot |_{B_{p,q}^\alpha(\nu)}$ is given by

$$|f|_{B_{p,q}^\alpha(\nu)} := \left(\int_0^\infty (t^{-\alpha} \omega_{s,L_p(\nu)}(f, t))^q \frac{dt}{t} \right)^{\frac{1}{q}} \quad , \quad 1 \leq q < \infty \quad ,$$

or

$$|f|_{B_{p,\infty}^\alpha(\nu)} := \sup_{t>0} (t^{-\alpha} \omega_{s,L_p(\nu)}(f, t)) \quad ,$$

see e.g. (Adams and Fournier, 2003, Section 7) and (Triebel, 2010, Sections 2 and 3). Note that $\|f\|_{B_{p,q}^\alpha(\nu)} := \|f\|_{L_p(\nu)} + |f|_{B_{p,q}^\alpha(\nu)}$ actually describes a norm of $B_{p,q}^\alpha(\nu)$ for all $q \in [1, \infty]$, see e.g. (DeVore and Lorentz, 1993, pp. 54/55) and (DeVore and Popov, 1988, p. 398). Again, if ν is the Lebesgue measure on X , we write $B_{p,q}^\alpha(X) := B_{p,q}^\alpha(\nu)$. For the sake of completeness, recall from e.g. (Adams and Fournier, 2003, Section 3) and (Triebel, 2010, Sections 2 and 3) the scale of Sobolev spaces $W_p^\alpha(\nu)$ defined by

$$W_p^\alpha(\nu) := \left\{ f \in L_p(\nu) : \partial^{(\beta)} f \in L_p(\nu) \text{ exists for all } \beta \in \mathbb{N}_0^d \text{ with } |\beta| \leq \alpha \right\},$$

where $\alpha \in \mathbb{N}_0$, $1 \leq p \leq \infty$, ν is an arbitrary measure, and $\partial^{(\beta)}$ is the β -th weak derivative for a multi-index $\beta = (\beta_1, \dots, \beta_d) \in \mathbb{N}_0^d$ with $|\beta| = \sum_{i=1}^d \beta_i$. That is, $W_p^\alpha(\nu)$ is the space of all functions in $L_p(\nu)$, whose weak derivatives up to order α exist and are contained in $L_p(\nu)$. Moreover, the Sobolev space is equipped with the Sobolev norm

$$\|f\|_{W_p^\alpha(\nu)}^p := \sum_{|\beta| \leq \alpha} \left\| \partial^{(\beta)} f \right\|_{L_p(\nu)}^p,$$

(cf. Adams and Fournier, 2003, page 60). We write $W_p^0(\nu) = L_p(\nu)$ and, for the Lebesgue measure μ on $X \subset \mathbb{R}^d$, we define $W_p^\alpha(X) := W_p^\alpha(\mu)$. It is well-known, see e.g. (Edmunds and Triebel, 1996, p. 25 and p. 44), that the Sobolev spaces $W_p^\alpha(\mathbb{R}^d)$ fall into the scale of Besov spaces, namely

$$W_p^\alpha(\mathbb{R}^d) \subset B_{p,q}^\alpha(\mathbb{R}^d)$$

for $\alpha \in \mathbb{N}$, $p \in (1, \infty)$, and $\max\{p, 2\} \leq q \leq \infty$. Moreover, for $p = q = 2$ we actually have equality, that is $W_2^\alpha(\mathbb{R}^d) = B_{2,2}^\alpha(\mathbb{R}^d)$ with equivalent norms.

Based on the least squares loss and RKHSs using Gaussian kernels over the partition sets A_j , the subsequent theorem refines the oracle inequality stated in Theorem 5.

Theorem 7 *Let $Y := [-M, M]$ for $M > 0$, $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ be the least squares loss and P be a distribution on $\mathbb{R}^d \times Y$. We write $X := \text{supp } P_X$. Furthermore, let (A) and (G) be satisfied. In addition, for an arbitrary subset $T \subset X$, we assume (T). Moreover, let $f_{L,P}^* : \mathbb{R}^d \rightarrow \mathbb{R}$ be a Bayes decision function such that $f_{L,P}^* \in L_2(\mathbb{R}^d) \cap L_\infty(\mathbb{R}^d)$ as well as $f_{L,P}^* \in B_{2,\infty}^\alpha(P_{X|A_T})$ for some $\alpha \geq 1$. Then, for all $p \in (0, 1)$, $n \geq 1$, $\tau \geq 1$, $\gamma = (\gamma_1, \dots, \gamma_m) \in (0, r]^m$, and $\lambda = (\lambda_1, \dots, \lambda_m) > 0$, the VP-SVM given by (4) using $\hat{H}_{\gamma_1}(A_1), \dots, \hat{H}_{\gamma_m}(A_m)$, and the loss L_{J_T} satisfies*

$$\begin{aligned} & \sum_{j=1}^m \lambda_j \|f_{D_j, \lambda_j, \gamma_j}\|_{\hat{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_{J_T}, P}(\hat{f}_{D, \lambda, \gamma}) - \mathcal{R}_{L_{J_T}, P}^* \\ & \leq C_{M, \alpha, p} \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} + r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p n^{-1} + \tau n^{-1} \right) \end{aligned}$$

with probability P^n not less than $1 - e^{-\tau}$, where $C_{M, \alpha, p} > 0$ is a constant only depending on M , α , p , d , $\|f_{L,P}^*\|_{L_2(\mathbb{R}^d)}$, $\|f_{L,P}^*\|_{L_\infty(\mathbb{R}^d)}$, and $\|f_{L,P}^*\|_{B_{2,\infty}^\alpha(P_{X|A_T})}$.

Using this oracle inequality, we derive learning rates w.r.t. the loss L_{J_T} for the learning method described by (3) and (4) in the following theorem.

Theorem 8 *Let $\tau \geq 1$ be fixed and $\beta \geq \frac{2\alpha}{d} + 1$. Under the assumptions of Theorem 7 and with*

$$r_n = c_1 n^{-\frac{1}{\beta d}}, \quad (21)$$

$$\lambda_{n,j} = c_2 r_n^d n^{-1}, \quad (22)$$

$$\gamma_{n,j} = c_3 n^{-\frac{1}{2\alpha+d}}, \quad (23)$$

for every $j \in \{1, \dots, m_n\}$, we have, for all $n \geq 1$ and $\xi > 0$,

$$\mathcal{R}_{L_{J_T}, P}(\hat{f}_{D, \lambda_n, \gamma_n}) - \mathcal{R}_{L_{J_T}, P}^* \leq C \tau n^{-\frac{2\alpha}{2\alpha+d} + \xi}$$

with probability P^n not less than $1 - e^{-\tau}$, where $\lambda_n := (\lambda_{n,1}, \dots, \lambda_{n,m_n})$ as well as $\gamma_n := (\gamma_{n,1}, \dots, \gamma_{n,m_n})$ and C, c_1, c_2, c_3 are positive constants with $c_3 \leq c_1$.

In the latter theorem the condition $\beta \geq \frac{2\alpha}{d} + 1$ is required to ensure $\gamma_{n,j} \leq r_n$, $j = 1, \dots, m_n$, which in turn is a prerequisite arising from Theorem 6 and the used entropy estimate. Let us briefly examine the extreme case $\beta = \frac{2\alpha}{d} + 1$. Using $r_n \approx n^{-\frac{1}{\beta d}}$ and (7) leads to covering numbers of the form $m_n \approx n^{\frac{d}{2\alpha+d}}$ and computational costs of $\mathcal{O}(m_n (\frac{n}{m_n})^q) = \mathcal{O}(n^{\frac{2\alpha q + d}{2\alpha+d}})$ which is actually less than the computational cost of order n^q , $q \in [2, 3]$, of an usual SVM. Note that for increasing β the computational cost of a VP-SVM is increasing as well. However, for $\beta > \frac{2\alpha}{d} + 1$, $r_n \approx n^{-\frac{1}{\beta d}}$, and $m_n \approx n^{\frac{1}{\beta}}$, a VP-SVM has costs of $\mathcal{O}(n^{\frac{1+(\beta-1)q}{\beta}})$ which still is less than $\mathcal{O}(n^q)$.

Let us finally take a closer look at the VP-SVM given by (4) and the considerations related to (17), where $f_{D, \lambda} \in H = \bigoplus_{j=1}^m \hat{H}_j$ solves the minimization problem

$$f_{D, \lambda} = \arg \min_{f_1 \in \hat{H}_1, \dots, f_m \in \hat{H}_m} \sum_{j=1}^m \lambda_j \|f_j\|_{\hat{H}_j}^2 + \mathcal{R}_{L, D}(\sum_{j=1}^m f_j).$$

Choosing $\lambda_1 = \dots = \lambda_m$, the VP-SVM problem can be understood as ℓ_2 -multiple kernel learning (MKL) problem using the RKHSs $\hat{H}_1, \dots, \hat{H}_m$. Learning rates for MKL have been treated, for example, in (Suzuki, 2011) and (Kloft and Blanchard, 2012). Assuming $f_{L, P}^* \in H$, the learning rate achieved in (Suzuki, 2011) is $mn^{-\frac{1}{1+s}}$ for dense settings, where s is the so-called spectral decay coefficient. In addition, Kloft and Blanchard (2012) obtain essentially the same rates under these assumptions. Let us therefore briefly investigate the above rate of (Suzuki, 2011). For RKHSs that are continuously embedded in a Sobolev space $W_2^\alpha(X)$, we have $s = \frac{d}{2\alpha}$ such that the learning rate reduces to $mn^{-\frac{2\alpha}{2\alpha+d}}$. Note that this learning rate is m times the optimal learning rate $n^{-\frac{2\alpha}{2\alpha+d}}$, where the number $m = m_n$ of kernels may increase with the sample size n . In particular, if $m_n \rightarrow \infty$ polynomially, then the rates obtained in (Suzuki, 2011) become substantially worse than the optimal rate. In contrast, due to the special choice of the RKHSs, this is not the case for our VP-SVM problem, provided that m_n does not grow faster than $n^{1/\beta}$.

Note that the oracle inequalities and learning rates achieved in Theorems 7 and 8 require $f_{L,P}^* \in B_{2,\infty}^\alpha(P_X|\bigcup_{j \in J_T} A_j)$. However, for an increasing sample size n , the sets A_j shrink and the index set J_T , indicating every set A_j such that $A_j \cap T \neq \emptyset$ and $T \subset \bigcup_{j \in J_T} A_j$, increases. In particular, this also involves that the set $\bigcup_{j \in J_T} A_j$ covering T changes in tandem with n . Since this is very inconvenient and since it would be desirable to assume a certain level of smoothness of the target function on a fixed region for all $n \in \mathbb{N}$, we consider the set T enlarged by an δ -tube. To this end, for $\delta > 0$, we define $T^{+\delta}$ by

$$T^{+\delta} := \{x \in X : \exists t \in T \text{ such that } \|x - t\|_2 \leq \delta\}, \quad (24)$$

which implies $T \subset T^{+\delta} \subset X$, cf. Figure 5. Note that, for every $\delta > 0$, there exists an $n_\delta \in \mathbb{N}$ such that for every $n \geq n_\delta$ the union of all partition sets A_j , having at least one common point with T , is contained in $T^{+\delta}$, i.e.

$$\forall \delta > 0 \quad \exists n_\delta \in \mathbb{N} \quad \forall n \geq n_\delta \quad : \quad \bigcup_{j \in J_T} A_j \subset T^{+\delta}, \quad (25)$$

where $J_T := \{j \in \{1, \dots, m_n\} : A_j \cap T \neq \emptyset\}$. Collectively, this implies

$$T \subset \bigcup_{j \in J_T} A_j \subset T^{+\delta}$$

for all $n \geq n_\delta$. Furthermore, since every set A_j is contained in a ball with radius $r_n = cn^{-\frac{1}{\beta d}}$, the lowest sample size n_δ in (25) can be determined by choosing the smallest $n_\delta \in \mathbb{N}$ such that $\delta \geq 2r_{n_\delta}$ with r_{n_δ} as in (7), that is

$$n_\delta = \left\lceil \left(\frac{2c}{\delta} \right)^{\beta d} \right\rceil.$$

This leads to the following corollary where we present an oracle inequality and learning rates assuming the smoothness level α of the target function on a fixed region.

Corollary 9 *Let $Y := [-M, M]$ for $M > 0$, $L : Y \times \mathbb{R} \rightarrow [0, \infty)$ be the least squares loss, and P be a distribution on $\mathbb{R}^d \times Y$. We write $X := \text{supp } P_X$. Furthermore, let (A) and (G) be satisfied. In addition, for an arbitrary subset $T \subset X$, we assume (T). Moreover, let $f_{L,P}^* : \mathbb{R}^d \rightarrow \mathbb{R}$ be a Bayes decision function such that $f_{L,P}^* \in L_2(\mathbb{R}^d) \cap L_\infty(\mathbb{R}^d)$ as well as*

$$f_{L,P}^* \in B_{2,\infty}^\alpha(P_X|_{T^{+\delta}})$$

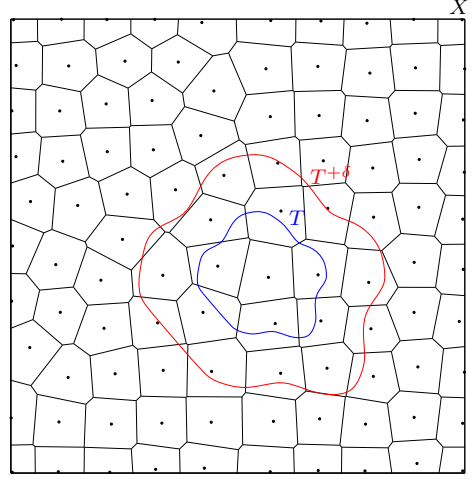


Figure 5: An input space X with the corresponding Voronoi partition as well as a subset $T \subset X$ enlarged by an δ -tube to $T^{+\delta}$.

for $\alpha \geq 1$ and some $\delta > 0$. Then, for all $p \in (0, 1)$, $n \geq n_\delta$, $\tau \geq 1$, $\gamma = (\gamma_1, \dots, \gamma_m) \in (0, r]^m$, and $\lambda = (\lambda_1, \dots, \lambda_m) > 0$, the VP-SVM given by (4) using $\hat{H}_{\gamma_1}(A_1), \dots, \hat{H}_{\gamma_m}(A_m)$, and the loss L_T satisfies

$$\begin{aligned} & \sum_{j=1}^m \lambda_j \|f_{D_j, \lambda_j, \gamma_j}\|_{\hat{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_T, P}(\hat{f}_{D, \lambda, \gamma}) - \mathcal{R}_{L_T, P}^* \\ & \leq C_{M, \alpha, p} \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} + r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p n^{-1} + \tau n^{-1} \right) \end{aligned}$$

with probability P^n not less than $1 - e^{-\tau}$, where $C_{M, \alpha, p} > 0$ is the same constant as in Theorem 7.

Additionally, let $\beta \geq \frac{2\alpha}{d} + 1$ as well as, for every $j \in \{1, \dots, m_n\}$, r_n , $\lambda_{n,j}$, and $\gamma_{n,j}$ be as in (21), (22), and (23), respectively, where c_1, c_2, c_3 are user-specified positive constants with $c_3 \leq c_1$. Then, for all $n \geq n_\delta = \left\lceil \left(\frac{2c_1}{\delta} \right)^{\beta d} \right\rceil$ and $\xi > 0$, we have

$$\mathcal{R}_{L_T, P}(\hat{f}_{D, \lambda_n, \gamma_n}) - \mathcal{R}_{L_T, P}^* \leq C \tau n^{-\frac{2\alpha}{2\alpha+d} + \xi}$$

with probability P^n not less than $1 - e^{-\tau}$, where $\lambda_n := (\lambda_{n,1}, \dots, \lambda_{n,m_n})$, $\gamma_n := (\gamma_{n,1}, \dots, \gamma_{n,m_n})$, and C is a positive constant.

Note that the assumption $f_{L,P}^* \in B_{2,\infty}^\alpha(P_{X|T+\delta})$ made in Corollary 9 is satisfied if, for example, $f_{L,P}^* \in B_{2,\infty}^\alpha(T^{+2\delta})$ and P_X has a bounded Lebesgue density on $T^{+\delta}$. Moreover, if this density is even bounded away from 0, it is well-known that the minmax rate is $n^{-\frac{2\alpha}{2\alpha+d}}$ for $\alpha > d/2$ and target functions $f_{L,P}^* \in W_2^\alpha(T)$. Modulo ξ , our rate is therefore asymptotically optimal in a minmax sense on T . In addition, for $\alpha > d$, the learning rates obtained for $f_{L,P}^* \in B_{2,\infty}^\alpha(T)$ are again asymptotically optimal modulo ξ on T .

5.2 Data-Dependent Parameter Selection for VP-SVMs

Note that in the previous theorems the choice of the regularization parameters $\lambda_{n,1}, \dots, \lambda_{n,m_n}$ and the kernel widths $\gamma_{n,1}, \dots, \gamma_{n,m_n}$ requires us to know the smoothness parameter α . Unfortunately, in practice, we usually do not know neither this value nor its existence. In this subsection, we thus show that a training/validation approach similar to the one examined in (Steinwart and Christmann, 2008a, Chapters 6.5, 7.4, 8.2) and (Eberts and Steinwart, 2013) achieves the same rates adaptively, i.e. without knowing α . For this purpose, let $\Lambda := (\Lambda_n)$ and $\Gamma := (\Gamma_n)$ be sequences of finite subsets $\Lambda_n \subset (0, r_n^d]$ and $\Gamma_n \subset (0, r_n]$. For a data set $D := ((x_1, y_1), \dots, (x_n, y_n))$, we define

$$\begin{aligned} D_1 &:= ((x_1, y_1), \dots, (x_l, y_l)), \\ D_2 &:= ((x_{l+1}, y_{l+1}), \dots, (x_n, y_n)), \end{aligned}$$

where $l := \lfloor \frac{n}{2} \rfloor + 1$ and $n \geq 4$. We further split these sets in data sets

$$D_j^{(1)} := \{(x_i, y_i) \in D_1 : x_i \in A_j\}, \quad j \in \{1, \dots, m_n\},$$

$$D_j^{(2)} := \{(x_i, y_i) \in D_2 : x_i \in A_j\}, \quad j \in \{1, \dots, m_n\},$$

and define $l_j := |D_j^{(1)}|$ for all $j \in \{1, \dots, m_n\}$ such that $\sum_{j=1}^{m_n} l_j = l$. For every $j \in \{1, \dots, m_n\}$, we basically use $D_j^{(1)}$ as a training set, i.e. based on D_1 in combination with the loss function $L_j := \mathbb{1}_{A_j} L$ we compute SVM decision functions

$$f_{D_j^{(1)}, \lambda_j, \gamma_j} := \arg \min_{f \in \hat{H}_{\gamma_j}(A_j)} \lambda_j \|f\|_{\hat{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_j, D_1}(f), \quad (\lambda_j, \gamma_j) \in \Lambda_n \times \Gamma_n.$$

Again, note that $f_{D_j^{(1)}, \lambda_j, \gamma_j} = 0$ if $D_j^{(1)} = \emptyset$. Next, for each j , we use D_2 in tandem with L_j (or essentially $D_j^{(2)}$) to determine a pair $(\lambda_{D_2, j}, \gamma_{D_2, j}) \in \Lambda_n \times \Gamma_n$ such that

$$\mathcal{R}_{L_j, D_2} \left(\hat{f}_{D_j^{(1)}, \lambda_{D_2, j}, \gamma_{D_2, j}} \right) = \min_{(\lambda_j, \gamma_j) \in \Lambda_n \times \Gamma_n} \mathcal{R}_{L_j, D_2} \left(\hat{f}_{D_j^{(1)}, \lambda_j, \gamma_j} \right).$$

Finally, combining the decision functions $f_{D_j^{(1)}, \lambda_{D_2, j}, \gamma_{D_2, j}}$ for all $j \in \{1, \dots, m_n\}$, and defining $\lambda_{D_2} := (\lambda_{D_2, 1}, \dots, \lambda_{D_2, m_n})$ and $\gamma_{D_2} := (\gamma_{D_2, 1}, \dots, \gamma_{D_2, m_n})$, we obtain a function

$$f_{D_1, \lambda_{D_2}, \gamma_{D_2}} := \sum_{j=1}^{m_n} f_{D_j^{(1)}, \lambda_{D_2, j}, \gamma_{D_2, j}} = \sum_{j=1}^{m_n} \mathbb{1}_{A_j} f_{D_j^{(1)}, \lambda_{D_2, j}, \gamma_{D_2, j}},$$

and we call every learning method that produces these resulting decision functions $f_{D_1, \lambda_{D_2}, \gamma_{D_2}}$ a *training validation Voronoi partition support vector machine* (TV-VP-SVM) w.r.t. $\Lambda \times \Gamma$. Moreover, using (5) we have, for $\lambda := (\lambda_1, \dots, \lambda_{m_n})$ and $\gamma := (\gamma_1, \dots, \gamma_{m_n})$,

$$\begin{aligned} \mathcal{R}_{L, D_2} \left(\hat{f}_{D_1, \lambda_{D_2}, \gamma_{D_2}} \right) &= \sum_{j=1}^{m_n} \mathcal{R}_{L_j, D_2} \left(\hat{f}_{D_j^{(1)}, \lambda_{D_2, j}, \gamma_{D_2, j}} \right) \\ &= \sum_{j=1}^{m_n} \min_{(\lambda_j, \gamma_j) \in \Lambda_n \times \Gamma_n} \mathcal{R}_{L_j, D_2} \left(\hat{f}_{D_j^{(1)}, \lambda_j, \gamma_j} \right) \\ &= \min_{(\lambda, \gamma) \in (\Lambda_n \times \Gamma_n)^{m_n}} \sum_{j=1}^{m_n} \mathcal{R}_{L_j, D_2} \left(\hat{f}_{D_j^{(1)}, \lambda_j, \gamma_j} \right) \\ &= \min_{(\lambda, \gamma) \in (\Lambda_n \times \Gamma_n)^{m_n}} \mathcal{R}_{L, D_2} \left(\hat{f}_{D_1, \lambda, \gamma} \right), \end{aligned}$$

where $f_{D_1, \lambda, \gamma} := \sum_{j=1}^{m_n} f_{D_j^{(1)}, \lambda_j, \gamma_j}$ with $(\lambda_j, \gamma_j) \in \Lambda_n \times \Gamma_n$ for all $j \in \{1, \dots, m_n\}$. In other words, the function $\hat{f}_{D_1, \lambda_{D_2}, \gamma_{D_2}}$ really minimizes the empirical risk \mathcal{R}_{L, D_2} w.r.t. the validation data set D_2 and the loss L , where the minimum is taken over all functions $\hat{f}_{D_1, \lambda, \gamma}$ with $(\lambda, \gamma) \in (\Lambda_n \times \Gamma_n)^{m_n}$.

The following theorem presents learning rates for the above described TV-VP-SVM.

Theorem 10 *Let $r_n := cn^{-\frac{1}{\beta d}}$ with constants $c > 0$ and $\beta > 1$. Under the assumptions of Theorem 7 we fix sequences $\Lambda := (\Lambda_n)$ and $\Gamma := (\Gamma_n)$ of finite subsets $\Lambda_n \subset (0, r_n^d]$*

and $\Gamma_n \subset (0, r_n]$ such that Λ_n is an $(r_n^d \varepsilon_n)$ -net of $(0, r_n^d]$ and Γ_n is a δ_n -net of $(0, r_n]$ with $\varepsilon_n \leq n^{-1}$ and $\delta_n \leq n^{-\frac{1}{2+d}}$. Furthermore, assume that the cardinalities $|\Lambda_n|$ and $|\Gamma_n|$ grow polynomially in n . Then, for all $\xi > 0$, $\tau \geq 1$, and $\alpha < \frac{\beta-1}{2}d$, the TV-VP-SVM producing the decision functions $f_{D_1, \lambda_{D_2}, \gamma_{D_2}}$ satisfies

$$\mathbb{P}^n \left(\mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda_{D_2}, \gamma_{D_2}}) - \mathcal{R}_{L_{J_T}, P}^* \leq c\tau n^{-\frac{2\alpha}{2\alpha+d} + \xi} \right) \geq 1 - e^{-\tau},$$

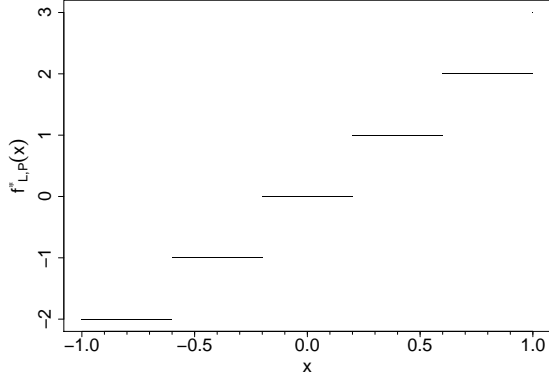
where $c > 0$ is a constant independent of n and τ .

Once more, we can replace the assumption $f_{L,P}^* \in B_{2,\infty}^\alpha(P_{X|A_T})$ by $f_{L,P}^* \in B_{2,\infty}^\alpha(P_{X|T+\delta})$ for some $\delta > 0$ and obtain the same learning rate as in Theorem 10 for all $n \geq n_\delta$ although $T^{+\delta}$ is fixed for all $n \in \mathbb{N}$. Note that, if P_X has a Lebesgue density that is bounded away from 0 and ∞ and either $f_{L,P}^* \in W_2^\alpha(T)$ for $\alpha > d/2$ or $f_{L,P}^* \in B_{2,\infty}^\alpha(T)$ for $\alpha > d$, these learning rates are again asymptotically optimal modulo ξ on T in a minmax sense. However, the condition $\alpha < \frac{\beta-1}{2}d$ restricts the set of α -values where we obtain learning rates adaptively. To be more precise, there is a trade-off between α and β . On the one hand, for small values of β only a small number of possible values for α is covered. On the other hand, for larger values of β the set of α -values where we achieve rates adaptively is increasing but the savings in terms of computing time is decreasing.

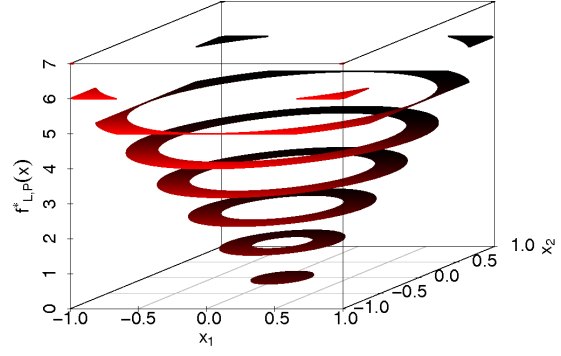
6. Experimental Results

In the previous sections we defined VP-SVMs and derived local learning rates that are essentially optimal. So far, it is, however, not clear if the theoretical results suggesting a generalization performance not worse than that of global SVMs can be empirically confirmed and if the predicted advantages of VP-SVMs in terms of computational costs are preserved in practice. Note that the latter is not as obvious as it may seem to be, since VP-SVMs create an overhead when generating the working sets, and the working sets themselves do not need to be as balanced as we assumed in our naïve analysis. In this section, we thus investigate the performance of VP-SVMs empirically. Namely, we carry out some experiments using the least squares loss with the objective to answer the subsequent questions:

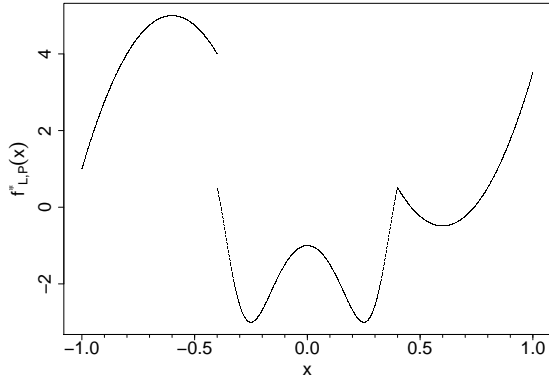
- (1) How do different radii affect the performance of VP-SVMs? In particular, what is the impact on the training time and the VP-SVM's test error?
- (2) How do the VP-LS-SVMs perform compared to the usual LS-SVMs in terms of the test error? What is the speed-up?
- (3) How does the performance of VP-SVMs compare to vanilla data splitting approaches such as random chunking (*RC-SVM*), in which the data set is divided into a random partition with equally sized subsets, and the final decision function is the average of the SVMs computed on each subset?
- (4) How does the VP-LS-SVM behave compared to the global LS-SVM, if the regression function has interruptions of its smoothness on zero sets?



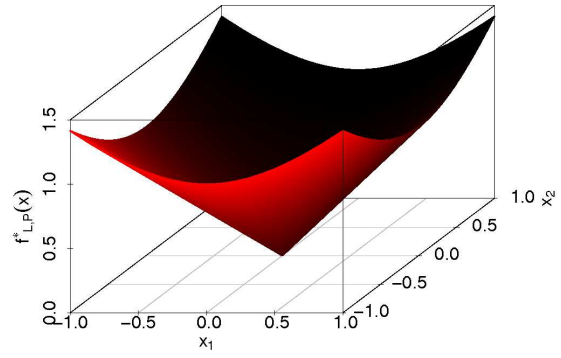
(a) Artificial data Type I: step function



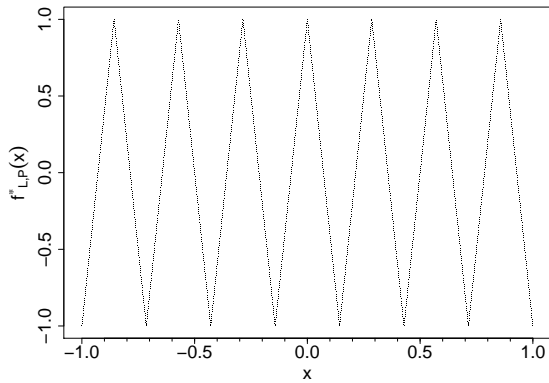
(d) Artificial data Type IV: circular step function



(b) Artificial data Type II: cracked function



(e) Artificial data Type V: 2-dimensional Euclidean norm



(c) Artificial data Type III: jagged function

Figure 6: Unscaled basic functions used to generate the artificial data sets.

To address these questions we utilize two kinds of data sets. On the one hand, to answer questions (1), (2), and (3), we examine the three real data sets COVTYPE, IJCNN1, and COD-RNA, which we obtained from LIBSVM’s homepage, see (Chang and Lin, 2011). Table 1 summarizes some characteristics of these data sets. On the other hand, we generated

data set type	full data set size	dimension	number of labels
COVTYPE	581 012	54	2
COD-RNA	488 565	8	2
IJCNN1	141 691	22	2

Table 1: Characteristics of the considered LIBSVM data sets.

several artificial data sets to address the last question. In order to prepare the data sets for the experiments, we edited the data sets from LIBSVM in the following manner. If for a real-world data set type the raw data set was already split, we first merged these sets so that we obtained one data set for each data set type. In a next step, we scaled the data componentwise such that all samples including labels lie in $[-1, 1]^{d+1}$, where d is the dimension of the input data. Finally, for each data set type, we generated random subsets that were afterwards randomly splitted into a training and a test data set. In this manner, we obtained, for each of the three LIBSVM data set types, training sets consisting of $n = 1\,000, 2\,500, 5\,000, 10\,000, 25\,000, 50\,000, 100\,000$ samples. Additionally, for the data sets COVTYPE and COD-RNA, we created training sets of sizes 250 000 and 500 000, and of sizes 250 000 and 400 000, respectively. The test data sets associated to the various training sets consist of $n_{\text{test}} = 50\,000$ random samples, apart from the training sets with $n_{\text{train}} \leq 5\,000$, for which we took $n_{\text{test}} = 10\,000$ test samples.

For the artificial data, we proceeded in a slightly different way. To generate the data sets we took as fundament the five regression functions pictured in Figure 6 and as noise, the sum of two uniform distributions on $[-c(x), c(x)]$, where $c(x) = \frac{1}{4} (3 \sin(\frac{\pi}{2}|x|) + 1)$ for the one-dimensional data sets and $c(x) = \frac{1}{4} (\sin(\frac{\pi}{4}(|x_1| + |x_2|)) + 1)$ for the two-dimensional data sets. Thus, we produced five different types of artificial data sets, where the various data set types are named according to their type numbers as in Figure 6. Initially, we created two sets, namely one training and one test data set, each consisting of 10 000 random input samples contained in $[-1, 1]$ and $[-1, 1]^2$, respectively. Then, for each artificial data set type, we determined the labels belonging to the input data as sum of the corresponding functional value and the noise and, finally, scaled all 20 000 labels to $[-1, 1]$. In a last step, we randomly built subsets of the training sets of size $n = 1\,000, 2\,500, 5\,000$. In this way, we altogether obtained, for each type of artificial data, four training data sets of size $n = 1\,000, 2\,500, 5\,000, 10\,000$ and a corresponding test data set of size $n_{\text{test}} = 10\,000$. Based on the test data sets the Bayes risks can be determined, see Table 2 where the Bayes risks are summarized for the various artificial data set types.

	Type I	Type II	Type III	Type IV	Type V
Bayes risk	0.0254	0.0137	0.0529	0.0083	0.0634

Table 2: Bayes risks w.r.t. test data sets for the various artificial data set types.

To minimize random effects, we repeated the experiment for each setting several times. Since experiments using large data sets entail long run times, we reran every experiment using a training set of size $n \geq 50\,000$ only three times while for training sets of size $n = 10\,000$, $25\,000$ we performed ten repetitions and for smaller training sets, namely of size $n = 1\,000$, $2\,500$, $5\,000$, even 100 runs. An exception are the experiments using artificial training sets of size $n = 10\,000$, where we realized 100 repetitions for the sake of uniformity.

To approach the above problems we used the least squares loss and Gaussian kernels for all experiments. We implemented an LS-SVM-solver in C++ similar to the one in (Steinwart et al., 2011). Around this solver, we then built the routines for the VP-SVM and the RC-SVM. The compilation of the three programmes was executed by LINUX's gcc. To produce comparable results in terms of run time, all real-world data experiments were realized by the same professional compute server¹ equipped with four INTEL XEON E7-4830 (2.13 GHz) 8-core processor, 256 GB RAM, and a 64 bit version of Debian GNU/Linux 6.0.7. In order that we can indeed compare their run time, we used eight cores to pre-compute the kernel matrix and to evaluate the final decision functions on the test set, and one core for the subsequent solver for every real data experiment. Since the artificial data sets consist of at most 10 000 samples we performed the according experiments by a computer equipped with one INTEL CORE i7-3770K (3.50 GHz) quad core processor, 16 GB RAM, and a 64 bit version of Debian GNU/Linux 6.0.7. For all artificial data experiments we used four cores to pre-compute the kernel matrix and to evaluate the final decision functions on the test set, and again one core for the solver. Even with pre-computed kernel matrices, our experiments on the real-world data altogether required almost 810 hours (approximately 34 days) for training and additionally almost 4 days for testing. Moreover, the experiments on the artificial data took nearly 43 hours for training and 168 minutes for testing. Without pre-computing the kernel matrices, e.g. by applying a standard caching approach, preliminary experiments suggested a multiplication of the training time, which would have rendered the experiments infeasible. Besides, our experiments will show that the available amount of RAM does not restrict the size of the training sets used by an VP-SVMs as severely as the ones used by LS-SVMs.

Let us quickly illustrate the routines of the VP- and the RC-SVM implemented around the LS-solver. For the VP-SVM, we first split the training set by Algorithm 1 in several working sets representing a Voronoi partition w.r.t. the user-specified radius. For this purpose, Algorithm 1 initially determines a cover of the input data applying the farthest first traversal algorithm, see (Dasgupta, 2008) and (Gonzalez, 1985) for more details. Note that this procedure induces working sets whose sizes may be considerably varying. In the case of an RC-SVM the working sets are created randomly, where their sizes are basically equal and the number of working sets is predefined by the user. Then, for the VP-SVM- as well as for the RC-SVM-algorithm the implemented LS-solver is applied on every working set. For each working set, we randomly split the respective training data set of size n_{train} in five folds to apply 5-fold cross-validation in order to deal with the hyper-parameters λ and γ taken from an 10 by 10 grid geometrically generated in $[0.001 \cdot n_{\text{train}}^{-1}, 0.1] \times [0.5 \cdot n_{\text{train}}^{-1/d}, 10]$.

1. On this occasion, we would like to thank the Institute for Applied Analysis and Numerical Simulation of the University of Stuttgart, who placed the above mentioned compute server at our disposal and, thus, enabled us to realize our experiments on large real-world data sets. In consequence, the overall time available for our experiments was limited.

Algorithm 1 Determine a Voronoi partition of the input data

Require: Input data set $D_X = \{x_1, \dots, x_n\}$ with sample size $n \in \mathbb{N}$ and some radius $r > 0$.

Ensure: Working sets indicating a Voronoi partition of D_X .

```

1: Pick an arbitrary  $z \in D_X$ 
2:  $Cover_1 \leftarrow z$ 
3:  $m \leftarrow 1$ 
4: while  $\max_{x \in D_X} \|x - Cover\|_2 > r$  do
5:    $z \leftarrow \arg \max_{x \in D_X} \|x - Cover\|_2$ 
6:    $m \leftarrow m + 1$ 
7:    $Cover_m \leftarrow z$ 
8:    $WorkingSet_m \leftarrow \emptyset$ 
9: end while
10: for  $i = 1$  to  $n$  do
11:    $k \leftarrow \arg \min_{j \in \{1, \dots, m\}} \|x_i - Cover_j\|_2$ 
12:    $WorkingSet_k \leftarrow WorkingSet_k \cup \{x_i\}$ 
13: end for
14: return  $WorkingSet_1, \dots, WorkingSet_m$ 

```

Finally, we obtain one decision function for each working set. To further process these decision functions the VP-SVM-algorithms picks exactly one decision function depending on the working set affiliation of the input value. On the contrary, the RC-SVM-algorithm simply takes the average of all the decision functions. Moreover, since we scaled the labels of all data sets to $[-1, 1]$, the computed decision functions are clipped at ± 1 . Altogether, note that the usual LS-SVM-algorithm can be interpreted as special case of both the VP-SVM- and the RC-SVM-algorithm using one working set.

The experimental results for the three real data sets are summarized in Tables 3 to 6. These tables as well as Tables 7 to 11, containing the results for the experiments on the artificial data sets, can be found in the Appendix. In addition to the average run times of the training and test phases, these tables reflect inter alia the average test errors of the empirical SVM solutions. Additionally, the L_2 -errors of the empirical SVM solutions, i.e. the value of

$$\sqrt{\frac{1}{n_{\text{test}}} \sum_{i=1}^{n_{\text{test}}} \left(\hat{f}_{D, \lambda, \gamma}(x_{\text{test}i}) - f_{L, P}^*(x_{\text{test}i}) \right)^2},$$

is determined for the artificial data sets. Moreover, note that some of the result tables are incomplete for very large real-world training data sets. In these cases, the kernel matrix, whose size depends on the training set size, did not fit into the RAM of the used computer and, thus, these experiments were left out.

6.1 Experiments on Real-World Data

In this subsection, we adress questions (1), (2), and (3) by examining the results for the real-world data sets COVTYPE, COD-RNA, and IJCNN1, which are composed in Figures 7–9 and Tables 3–6.

6.1.1 COMPARISON OF VP-SVMs USING DIFFERENT RADII

In the following, we focus on the VP-SVMs using four different radii for the various real-world data sets, where the experimental results are summarized in Tables 3–6 as well as in Subfigures (d)–(f) of Figures 7–9. Examining the achieved training times for each data set type, we observe that, for increasing training set sizes, the radius that leads to the shortest training time typically decreases. More precisely, for the real data sets with sample size $n_{\text{train}} > 10\,000$, the VP-SVMs using the smallest radius always train fastest, while for the data sets with $n_{\text{train}} \leq 10\,000$, we can not make a uniform statement. Clearly, this finding is not surprising, since an SVM for a small data set trains considerably faster than an SVM for a large data set, such that splitting the large data set and running an SVM for each of the small data sets may altogether still be faster. Recall additionally the considerations in terms of the computational cost made in Section 2.

Let us now consider the VP-SVM results in terms of the realized test errors. As expected, for the real-world data sets, the test errors achieved by VP-SVMs with fixed radii decrease with increasing sample size of the used data sets except twice. In addition, for the real data sets COVTYPE and COD-RNA, the test errors decrease for increasing radii, cf. Figures 7(f) and 8(f). Here, however, the test errors achieved for the various radii get close to each other with increasing training sample size. The same behavior of the test errors appears for the IJCNN1 data set, though, for $n_{\text{train}} \geq 5\,000$, both intermediate radii yield even smaller empirical risks than the largest radius, see Figure 9(f). In consequence, it is not straightforward to draw any conclusion on the relation between radius and test error. Nonetheless, we can say that VP-SVMs using small radii enjoy test errors that are never significantly larger and sometimes even smaller than those of VP-SVMs using the largest of the applied radii.

Besides, Tables 3 and 5 or Figures 7(e),(f) and 8(e),(f) contain an additional finding. For large data sets, namely for the COVTYPE data set of size 500 000 or for the COD-RNA data sets of size 250 000 and 400 000, the VP-SVMs with large radii did not yield any solution, since they failed due to the technical requirements caused by the used computer. More precisely, in these cases, there was at least one working set such that its kernel matrix did not fit into the RAM any more. Fortunately, the working sets of VP-SVMs using smaller radii were small enough such that we still received an outcome. What is more, these VP-SVMs yielded a better empirical risk in partially less training time compared to VP-SVMs with large radii and training sample sizes that still allowed a successful performance. That is, using a small radius for the VP-SVM and a training set that is oversized for VP-SVMs with a larger radius reduces the test error. More precisely, a large training set is crucial for a small empirical risk, where the possibly arising computational restrictions can be eluded by a VP-SVM with an appropriate radius.

All in all, localized SVMs using some small radius lead in substantially less training time to either negligible worse or even better test errors than VP-SVMs with large radii, if the training sample size is adequate, i.e. $n_{\text{train}} \geq 5\,000$. In addition, the real data sets with a large sample size demonstrate that VP-SVMs with small radii are able to conquer the technical restrictions caused by the used computer and thus yield a better empirical risk than VP-SVMs with bigger radii can attain at all.

6.1.2 COMPARING VP-SVMs WITH GLOBAL LS-SVMs

In the following, we compare the results of the VP-SVM using different radii to the standard LS-SVM. For the real-world data sets COD-RNA and IJCNN1, the VP-SVMs, based on the largest of the applied radii, use only one working set. Thus, they coincide with the standard LS-SVM modulo different values generated by the random number generator. To verify this fact, we compare for the real data sets COD-RNA and IJCNN1 the results of the VP-SVMs using one working set to the results of the standard LS-SVMs, see Tables 5 and 6. Here, we note that the LS-SVM test errors typically decrease with increasing training sample size. The same holds for the VP-SVMs using one working set. Moreover, the latter VP-SVM and the LS-SVMs perform equally well in terms of training time and empirical risk, however, for $n_{\text{train}} \geq 25\,000$ the VP-SVMs train slower.

In practice, a crucial problem is caused by the run time required by an algorithm. Hence, for each data set type, we compare hereafter the LS-SVM to the VP-SVM that trains fastest for the largest training data set. The required average training times and the average test errors of these SVMs are illustrated in Subfigures (g)–(i) of Figures 7–9. First, we notice that the selected VP-SVM uses the smallest of the applied radii for each data set type. Besides, the LS-SVM’s test errors are lower than those of the VP-SVMs. However, with increasing training set size the VP-SVM’s test errors get close to the ones of the LS-SVM. Moreover, for the IJCNN1 data set of size 100 000, their empirical risks even coincide, cf. Figure 9(i). Besides, the VP-SVMs train considerably faster than the LS-SVMs. In particular, for $n_{\text{train}} = 100\,000$ the VP-SVMs require at most 8.5% of the LS-SVM’s training times, see Figures 7(h), 8(h), and 9(h). Finally, recall that, for data sets of size $n_{\text{train}} \geq 250\,000$, the LS-SVM problem is infeasible with our computer, just like the VP-SVMs using the largest of the applied radii. In contrast, for $n_{\text{train}} \geq 250\,000$, VP-SVMs using small radii usually train considerably faster and achieve lower test errors than the LS-SVMs for $n_{\text{train}} = 100\,000$, cf. Figures 7(h)–(i) and 8(h)–(i).

Concluding, we have seen that the application of a VP-SVM using a small radius instead of the standard LS-SVM reduces the run time considerably entailing at most a negligible worsening or even an improvement of the test errors. Moreover, applying VP-SVMs with sufficiently small radii enables us to use large data sets and, thus, to elude the computational restrictions to sufficiently small data sets. As a result, handling really large data sets with the help of suitable VP-SVMs can lead to significantly improved test errors compared to an LS-SVM setting with memory constraints.

6.1.3 COMPARISON OF VP-SVMs WITH RC-SVMs

First of all, let us investigate the RC-SVM results that are composed in Tables 4–6 as well as in Subfigures (a)–(c) of Figures 7–9. For the real data sets COVTYPE we considered ten, for the data sets COD-RNA nine, and for the data sets IJCNN1 eight different numbers of working sets. In each case, we started with an RC-SVM using one working set, i.e. with an RC-SVM that corresponds to the global LS-SVM modulo different values generated by the random generator, cf. Tables 5 and 6. Comparing for every data set the RC-SVMs using various numbers of working sets, we observe that the number of working sets, minimizing the RC-SVMs training time, increases in tandem with the sample size. Moreover, the RC-SVM using one working set never trains fastest compared to the other RC-SVMs using

more than one working set. Furthermore, the average test errors for the applied RC-SVMs usually decrease for a decreasing number of working sets and, hence, are minimized by the smallest possible number of working sets. Of course, all these findings are not surprising, since RC-SVMs are typically used to reduce the training time.

Let us now compare the results of VP- and RC-SVMs using roughly the same number of working sets, cf. Tables 3–6. Initially note that, even though we consider VP- and RC-SVM based on the same number of working sets, the RC-SVM working sets are about the same size whereas the VP-SVM working sets may have different sizes with a large range. That is, the VP-SVMs often deal with a few substantially larger working sets than the RC-SVMs. Consequently, the RC-SVMs often perform faster than the VP-SVMs, which require up to five times the RC-SVM’s training time for $n_{\text{train}} = 100\,000$. Contrarily, the average empirical risks achieved by the VP-SVMs are substantially lower than those of the RC-SVMs. Besides, in a few cases the VP-SVMs possess at least one working set which is oversized for the computer’s RAM, so that these VP-SVM problems are infeasible, whereas the comparable RC-SVMs avoid this conflict. Here, consider e.g. the RC-SVM using seven working sets and the VP-SVM with radius $r = 4$ for the COVTYPE data set of size 500 000.

In Section 6.1.2, we compared for each data set type the LS-SVM with the VP-SVM that trains fastest for the largest training data set. Here, we additionally compare this VP-SVM to the RC-SVMs. To be able to draw a fair comparison in terms of the achieved test errors, we choose those RC-SVMs that train roughly as fast as the VP-SVM for the largest training set, i.e. the slowest RC-SVM training faster and the fastest RC-SVM training slower than the above VP-SVM. Subfigures (g)–(i) of Figures 7–9 illustrate the average training times and the average test errors of these RC-SVMs, the above VP-SVM, and the LS-SVM. Considering the RC-SVMs, the faster of the two requires for $n_{\text{train}} = 100\,000$ between 51% and 83% of the VP-SVMs training time and trains at most seven minutes faster than the VP-SVM. However, at least for $n_{\text{train}} \geq 5\,000$, both considered RC-SVMs induce substantially higher test errors than the VP- and LS-SVM. Finally, note that VP-SVMs for $n_{\text{train}} \geq 250\,000$ considerably outperform LS-SVMs for $n_{\text{train}} = 100\,000$, while RC-SVMs for $n_{\text{train}} \geq 250\,000$ lead to even worse test errors than the considered LS-SVMs.

Summarizing, we record that RC-SVMs using as few as possible working sets achieve the smallest RC-SVM test errors, however, those using more working sets perform faster. Furthermore, compared to VP-SVMs using roughly the same number of working sets as the RC-SVMs, the latter ones may learn faster though not as good as the VP-SVMs. Moreover, considering RC-SVMs that require roughly the same training time as the fastest VP-SVM, we saw that the RC-SVMs lead to much higher empirical risks. That is, if the required training time is a hard constraint, then the VP-SVM that satisfies this constraint achieves a better test error than a RC-SVM that also trains fast enough.

6.2 Experiments on Artificial Data

It remains to address the last question. To this end, we consider the results on the various artificial data sets, on the one hand, for the LS-SVM and, on the other hand, for the VP-SVM performing fastest for $n_{\text{train}} = 10\,000$. Moreover, for the sake of comparability, we again add to this selection the two RC-SVMs training roughly as fast as the VP-SVM for the artificial data sets of size $n_{\text{train}} = 10\,000$. However, for the artificial data sets of Type I,

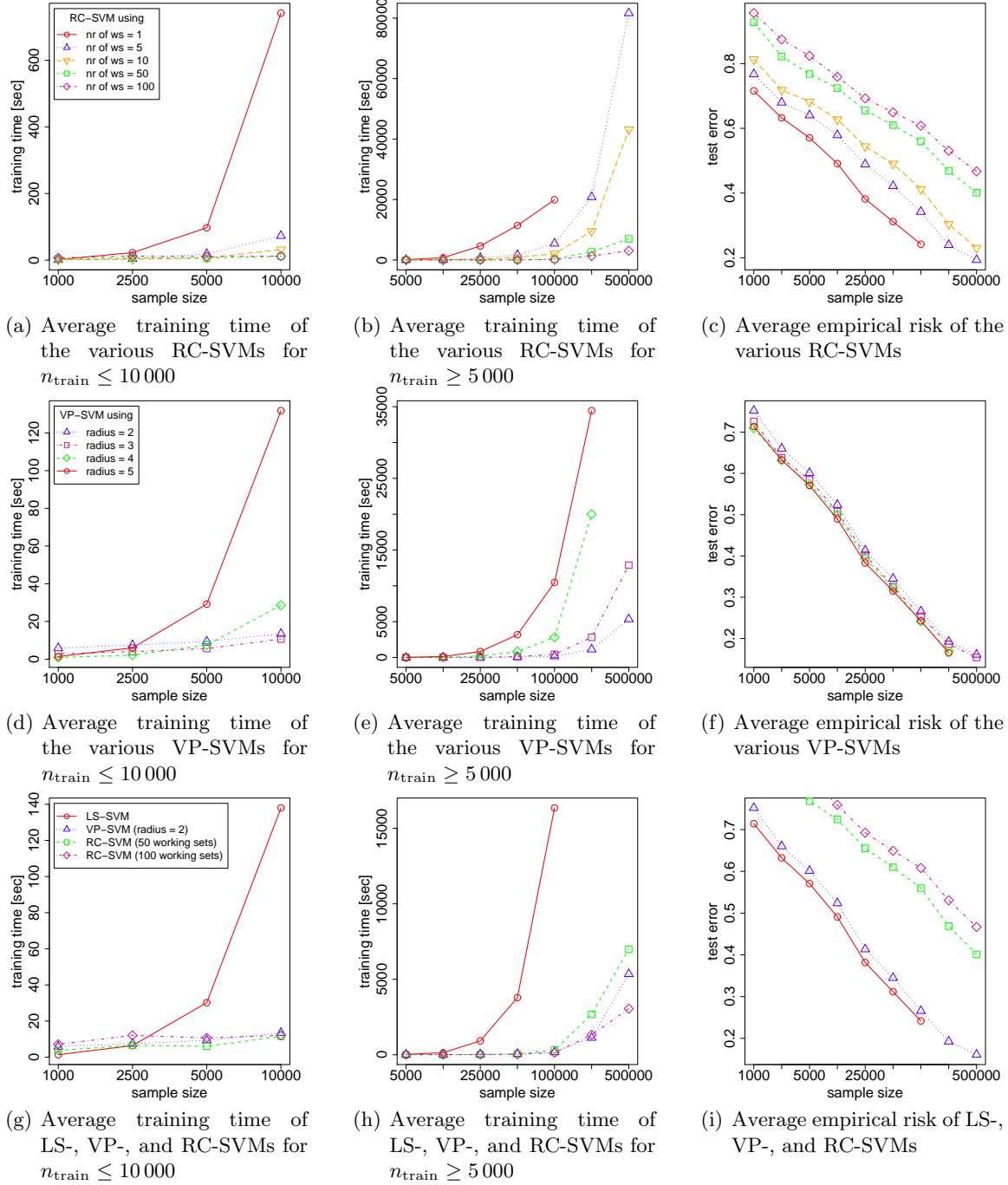


Figure 7: Average training time and test error of LS-, VP-, and RC-SVMs for the real-world data COVTYPE depending on the training set size $n_{\text{train}} = 1000, \dots, 500\,000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and two RC-SVMs. Here, the VP-SVM is the one which trains fastest for $n_{\text{train}} = 500\,000$ and the two RC-SVMs are those which achieve for $n_{\text{train}} = 500\,000$ roughly the same training time as the chosen VP-SVM. Here, note that, for $n_{\text{train}} = 10\,000$, the RC-SVM using one working set trains substantially slower than the LS-SVM, even though this RC-SVM is basically an LS-SVM. As a reason for this phenomenon, we conjecture that the used compute server was busy because of other influences.

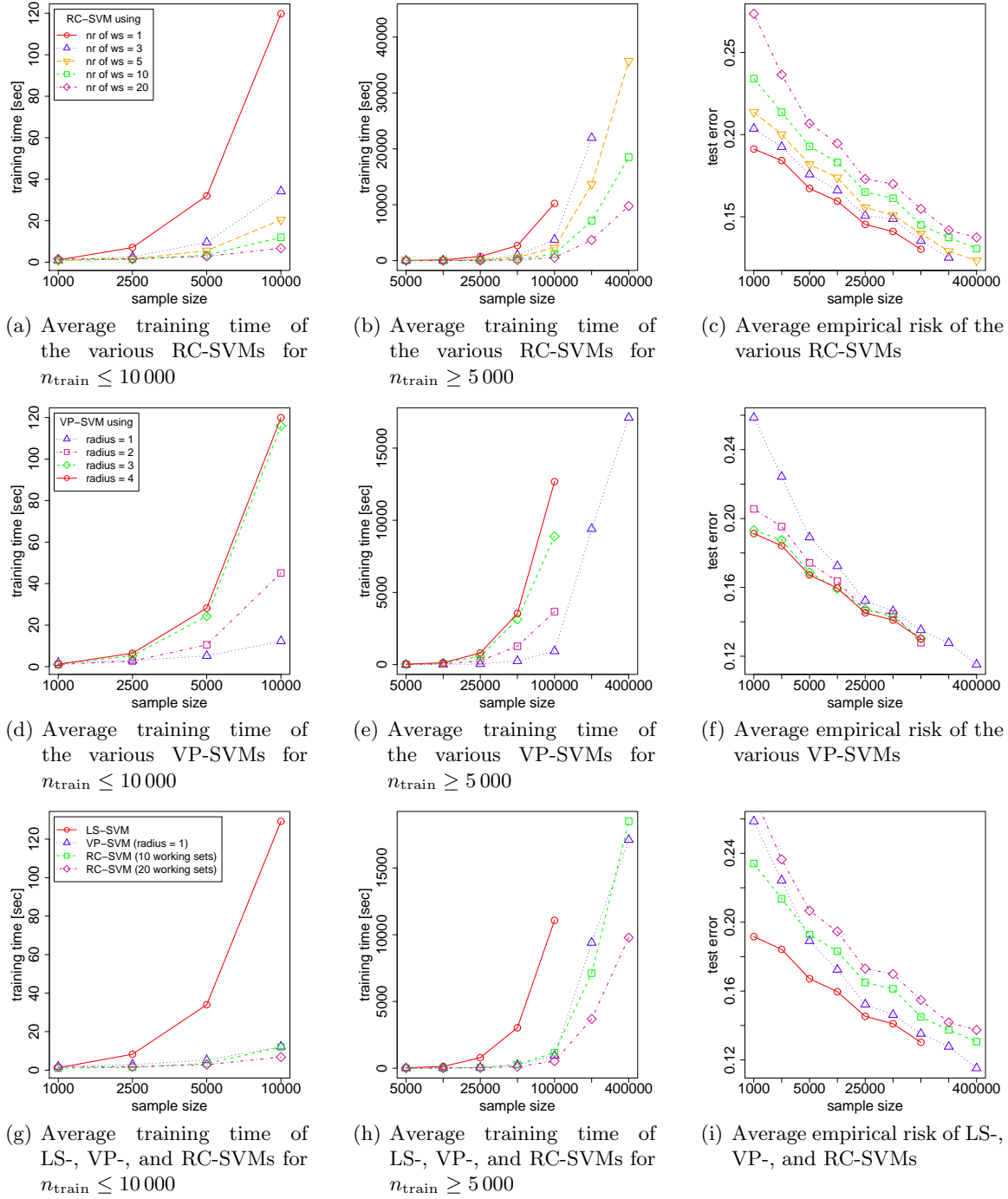


Figure 8: Average training time and test error of LS-, VP-, and RC-SVMs for the real-world data COD-RNA depending on the training set size $n_{\text{train}} = 1\,000, \dots, 400\,000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and two RC-SVMs. Here, the VP-SVM is the one which trains fastest for $n_{\text{train}} = 400\,000$ and the two RC-SVMs are those which achieve for $n_{\text{train}} = 400\,000$ roughly the same training time as the chosen VP-SVM.

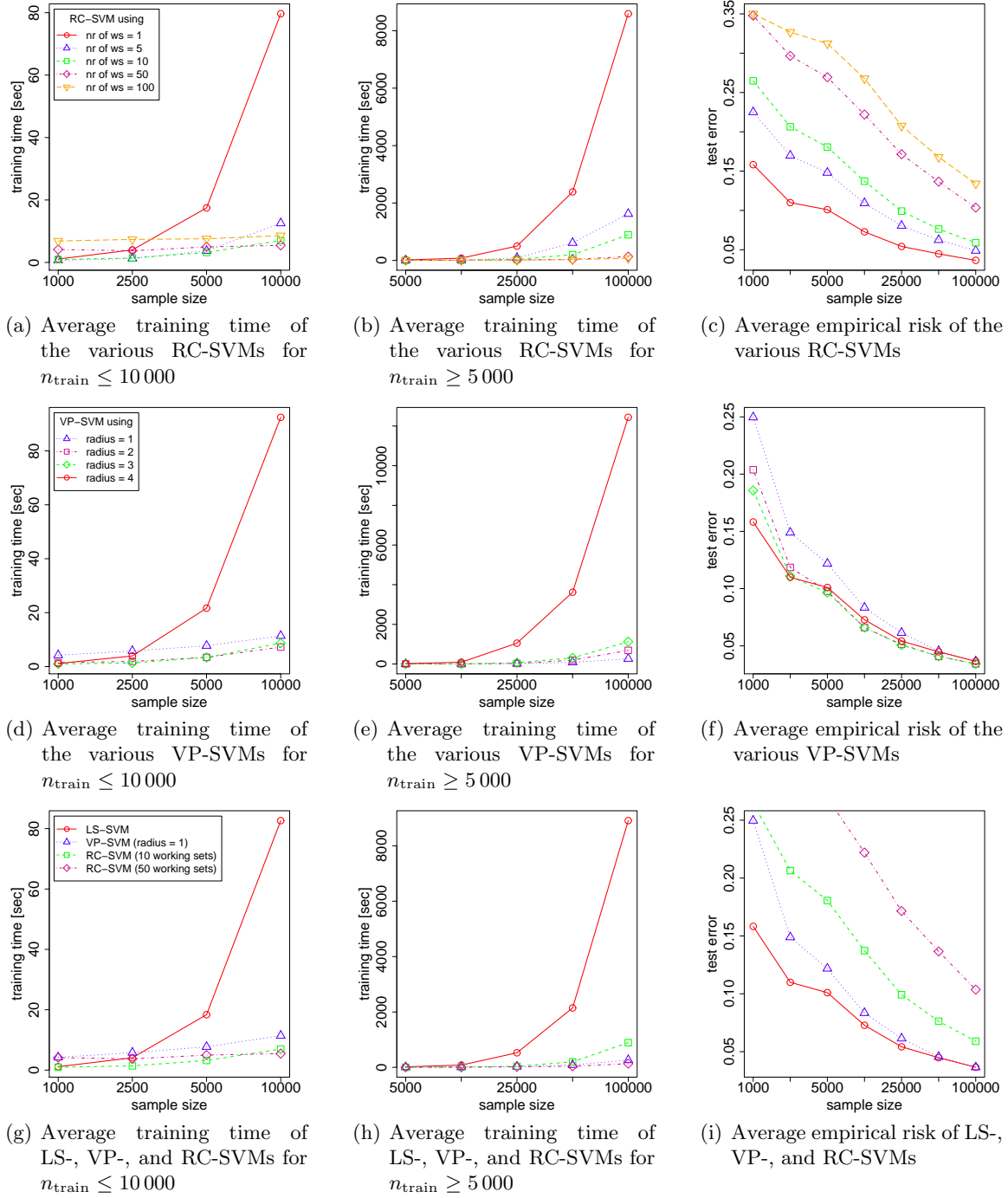


Figure 9: Average training time and test error of LS-, VP-, and RC-SVMs for the real-world data IJCNN1 depending on the training set size $n_{\text{train}} = 1\,000, \dots, 100\,000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and two RC-SVMs. Here, the VP-SVM is the one which trains fastest for $n_{\text{train}} = 100\,000$ and the two RC-SVMs are those which achieve for $n_{\text{train}} = 100\,000$ roughly the same training time as the chosen VP-SVM.

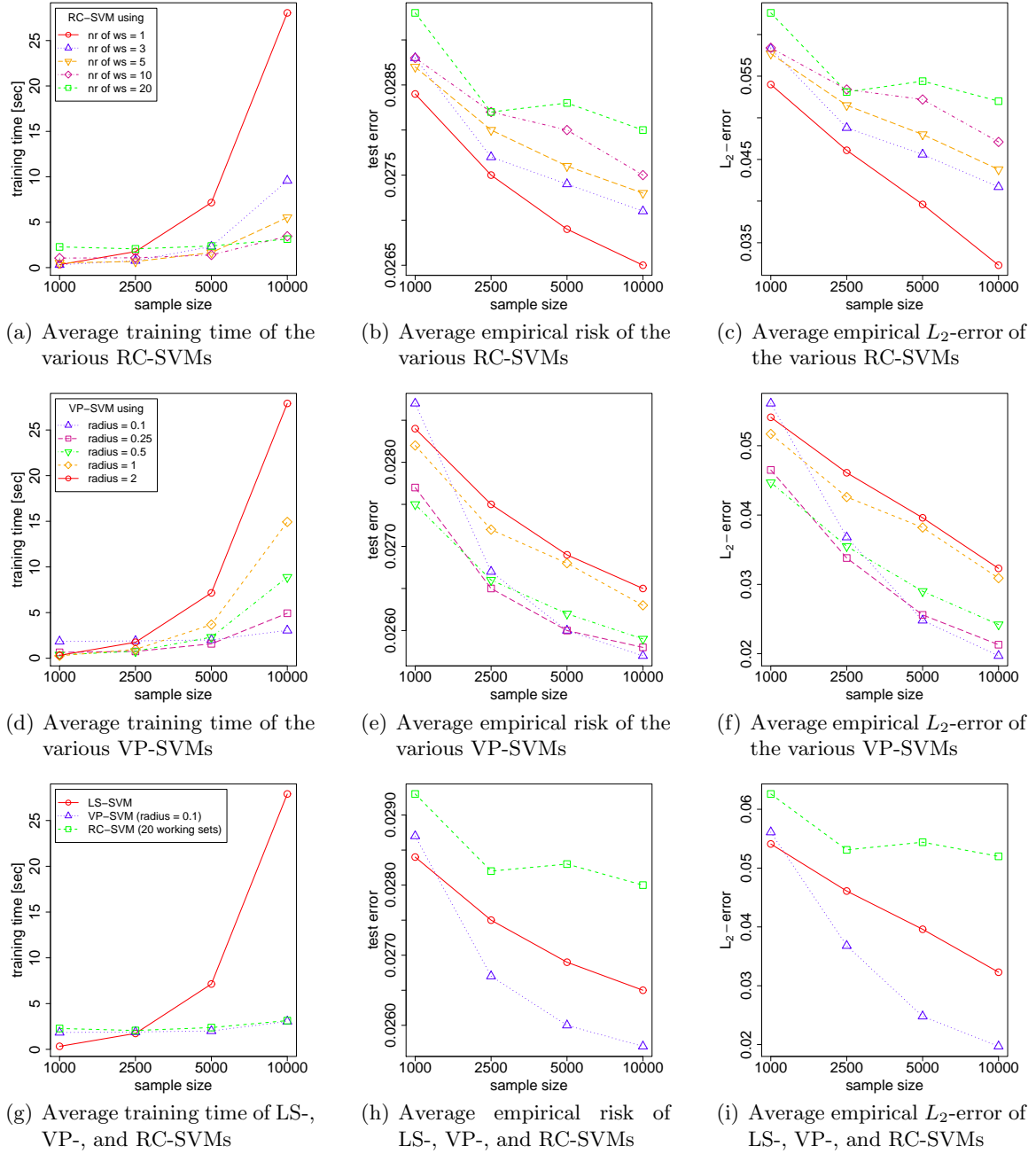


Figure 10: Average training time and test error of LS-, VP-, and RC-SVMs for the artificial data Type I depending on the training set size $n_{\text{train}} = 1000, \dots, 10000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and one RC-SVM. Here, the VP-SVM and the RC-SVM are those which train fastest for $n_{\text{train}} = 10000$. Note that in the case at hand none of the considered RC-SVMs performs faster than the fastest VP-SVM for $n_{\text{train}} = 10000$.

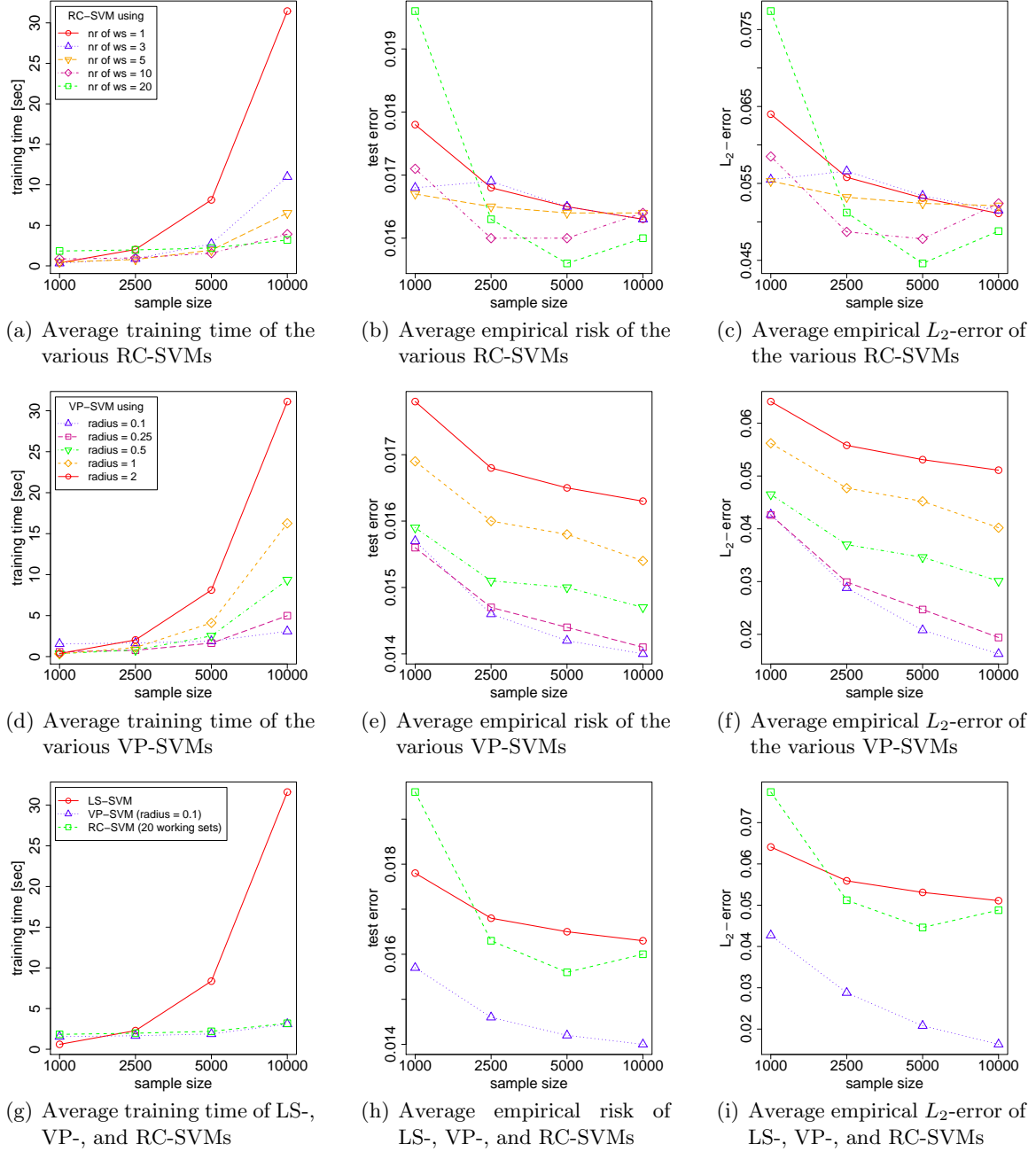


Figure 11: Average training time and test error of LS-, VP-, and RC-SVMs for the artificial data Type II depending on the training set size $n_{\text{train}} = 1000, \dots, 10000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and one RC-SVM. Here, the VP-SVM and the RC-SVM are those which train fastest for $n_{\text{train}} = 10000$. Note that in the case at hand none of the considered RC-SVMs performs faster than the fastest VP-SVM for $n_{\text{train}} = 10000$.

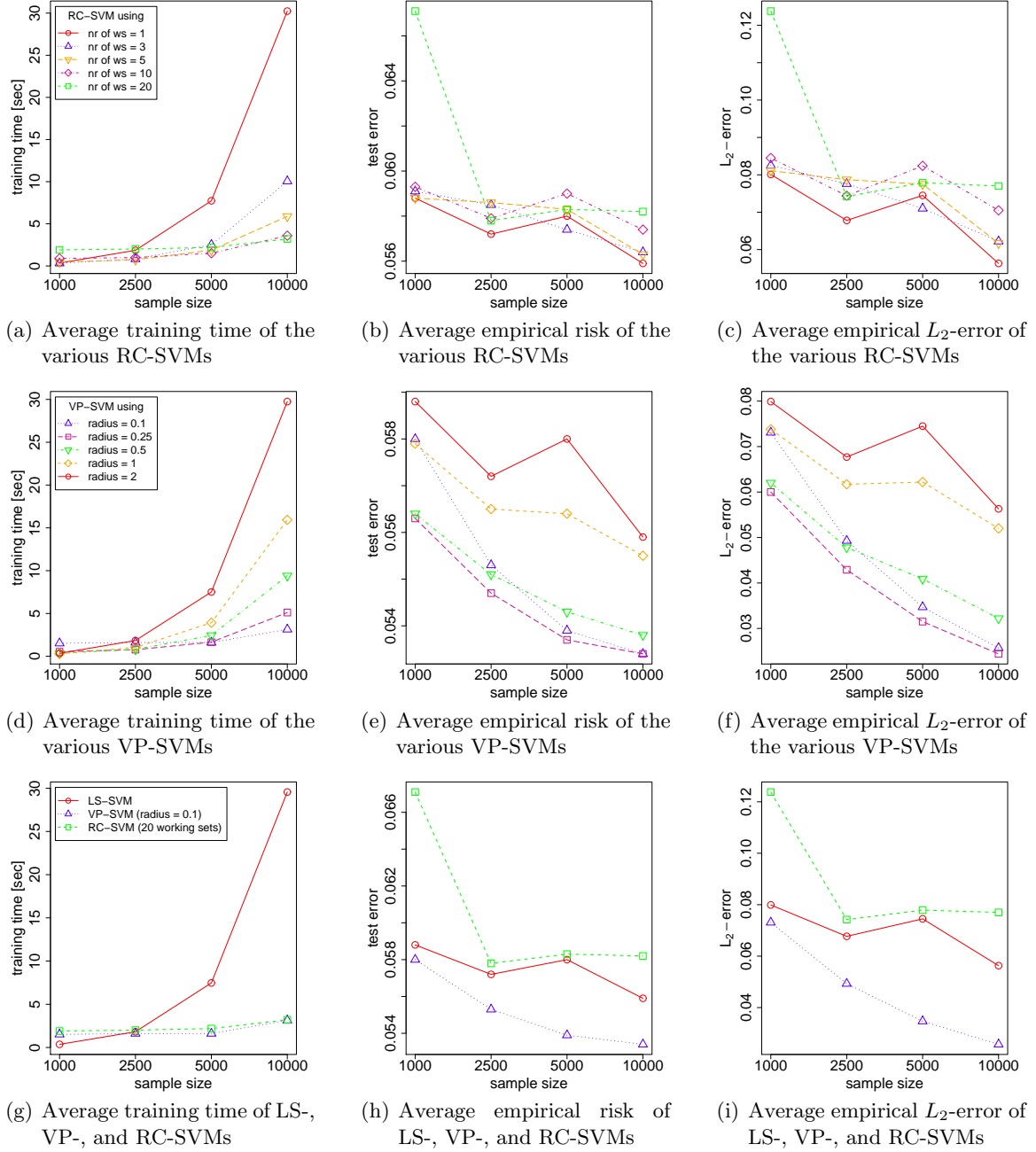


Figure 12: Average training time and test error of LS-, VP-, and RC-SVMs for the artificial data Type III depending on the training set size $n_{\text{train}} = 1000, \dots, 10000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and one RC-SVM. Here, the VP-SVM and the RC-SVM are those which train fastest for $n_{\text{train}} = 10000$. Note that in the case at hand none of the considered RC-SVMs performs faster than the fastest VP-SVM for $n_{\text{train}} = 10000$.

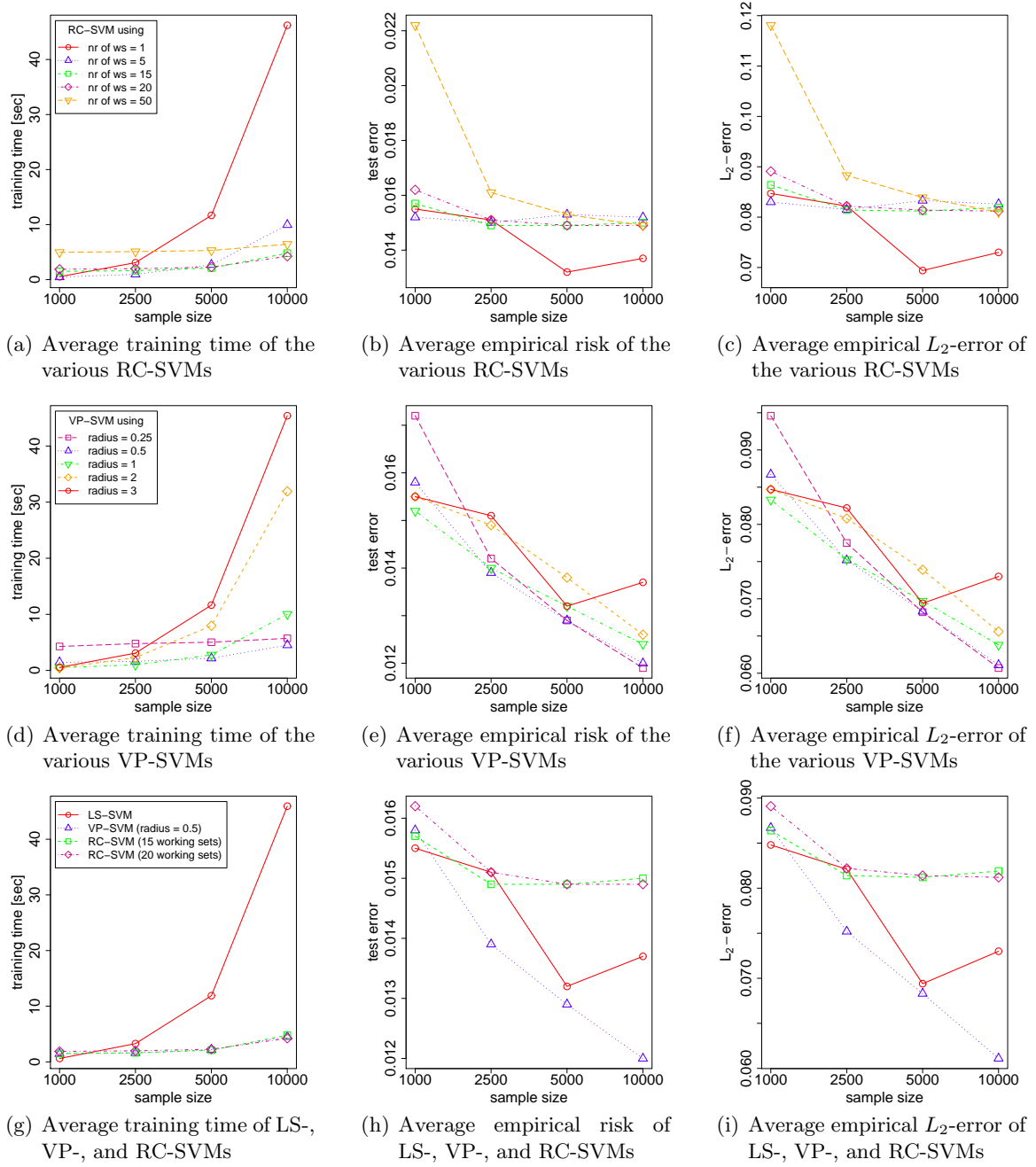


Figure 13: Average training time and test error of LS-, VP-, and RC-SVMs for the artificial data Type IV depending on the training set size $n_{\text{train}} = 1000, \dots, 10000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and two RC-SVMs. Here, the VP-SVM is the one which trains fastest for $n_{\text{train}} = 10000$ and the two RC-SVMs are those which achieve for $n_{\text{train}} = 10000$ roughly the same training time as the chosen VP-SVM.

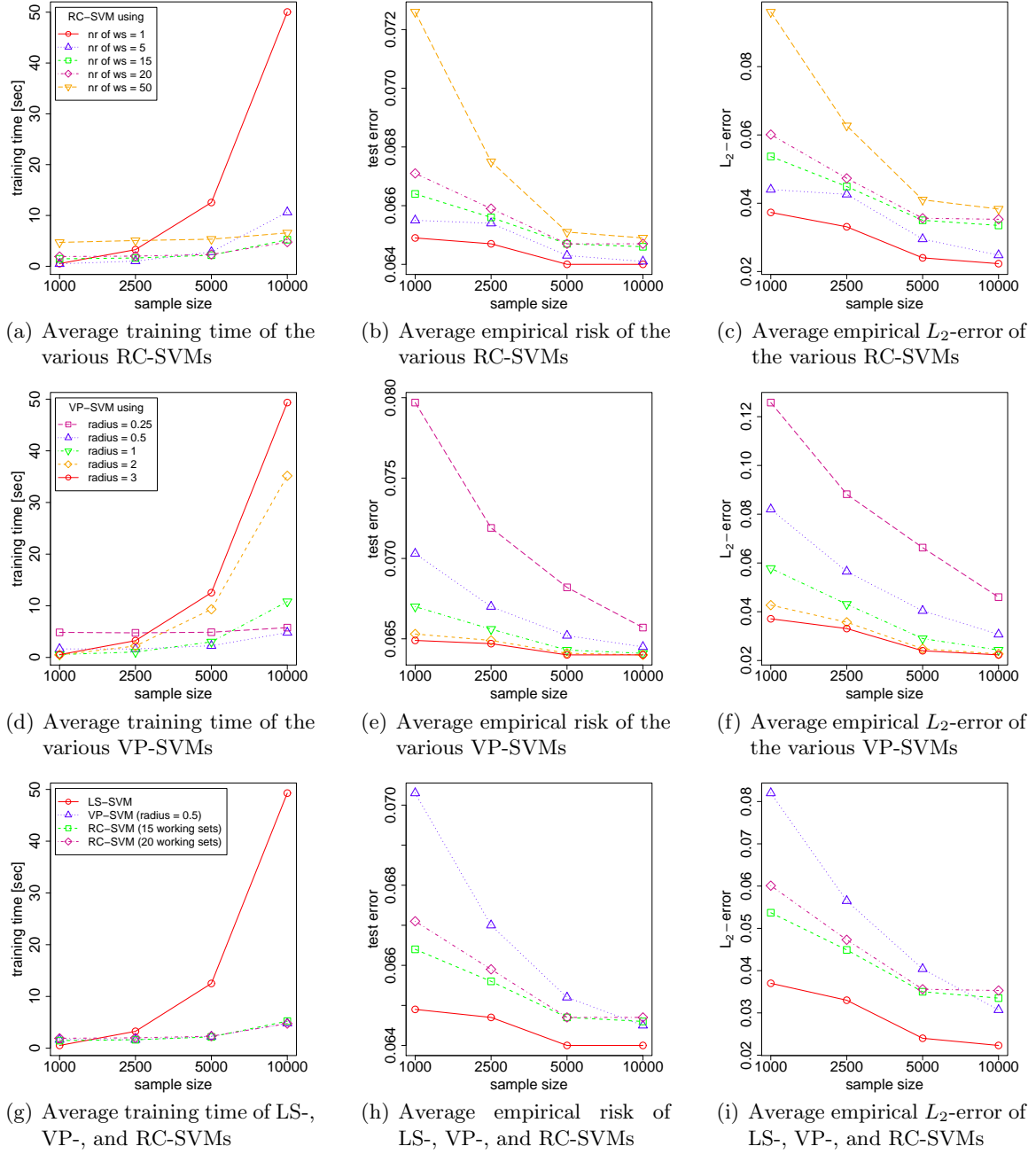


Figure 14: Average training time and test error of LS-, VP-, and RC-SVMs for the artificial data Type V depending on the training set size $n_{\text{train}} = 1000, \dots, 10000$. Subfigures (a)–(c) show the results for RC-SVMs using different numbers of working sets and Subfigures (d)–(f) illustrate the results for VP-SVMs using various radii. At the bottom, Subfigures (g)–(i) contain the average training times and the average test errors of the LS-SVM, one VP-SVM and two RC-SVMs. Here, the VP-SVM is the one which trains fastest for $n_{\text{train}} = 10000$ and the two RC-SVMs are those which achieve for $n_{\text{train}} = 10000$ roughly the same training time as the chosen VP-SVM.

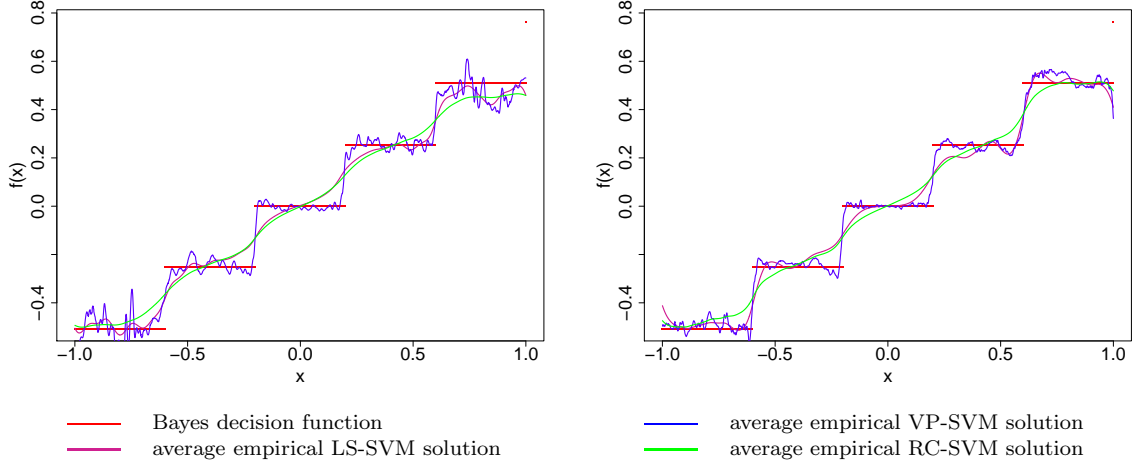


Figure 15: Predictions for the artificial data sets of Type I, drawn from the step function in Figure 6(a) with noise depending on x . The left graphic shows the predictions for the data set of size $n_{\text{train}} = 1000$ and the right graphic for the data set of size $n_{\text{train}} = 10000$. Here, note that the VP-SVM solutions are not necessarily continuous, nevertheless we continuously connected its predicted values in the above plots.

II, and III, none of the executed RC-SVMs trained faster for $n_{\text{train}} = 10000$ than the VP-SVM with the smallest radius, so that we only consider one RC-SVM in these cases. The required average training times and average test errors of the selected SVMs are illustrated in Subfigures (g)–(i) of Figures 10–14 and summarized in Tables 7–11. Here, we note that, for the artificial data of Type I, II, and III, the VP-SVM using the smallest of the applied radii trains fastest for $n_{\text{train}} = 10000$, while for the artificial data of Type IV and V it is the VP-SVM using the second smallest radius.

Expectedly, we detect an evident improvement of the various average empirical SVM solutions using 10000 training samples instead of 1000 samples. Besides, the considered VP-SVM trains substantially faster than the standard LS-SVM with less than 11% of the LS-SVM’s training time for $n_{\text{train}} = 10000$. Additionally, the VP-SVM’s test errors are usually considerably lower than the test errors of the LS-SVM. Regarding the test errors of the RC-SVMs, we note that, in the majority of cases, they are higher than the VP-SVM’s and the LS-SVM’s test errors.

So far, we examined the behavior of LS-, VP-, and RC-SVMs in terms of training time and test error. Let us finally compare the three different kinds of SVMs w.r.t. their optical appearance. To this end, the average empirical SVM solutions are plotted in Figures 15–20 for the different artificial data sets of size $n_{\text{train}} = 1000$ and 10000. Here, note that, for the artificial data of Type IV and V, we do not consider both RC-SVMs training roughly as fast as the selected VP-SVM but only the one of the both RC-SVMs with the lower test error.

The observation that, for the artificial data of Type I, II, and III, the VP-SVMs perform best, is reinforced by the average empirical VP-SVM solutions illustrated in Figures 15–18. More precisely, Figure 15 shows that only the VP-SVMs exhaust the widths of the steps

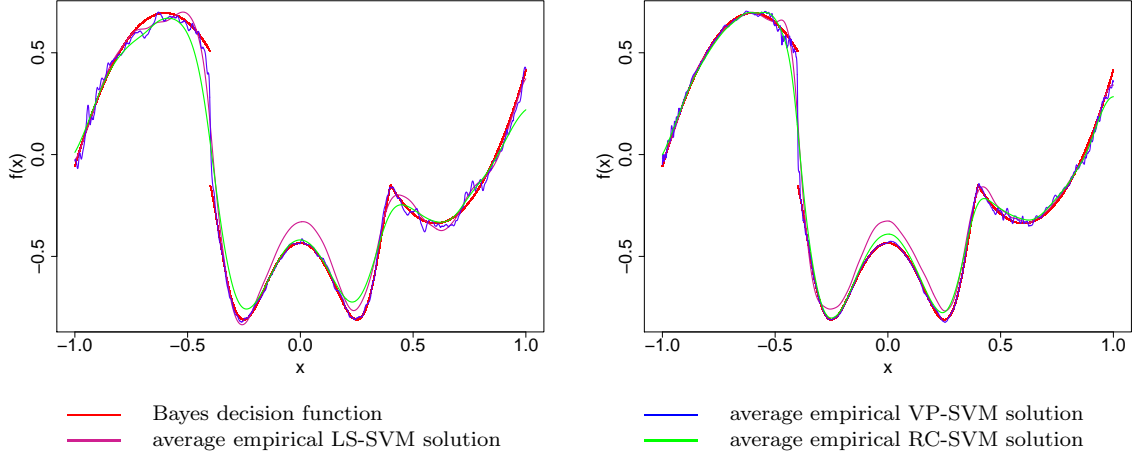


Figure 16: Predictions for the artificial data sets of Type II, drawn from the cracked function in Figure 6(b) with noise depending on x . The left graphic shows the predictions for the data set of size $n_{\text{train}} = 1000$ and the right graphic for the data set of size $n_{\text{train}} = 10000$.

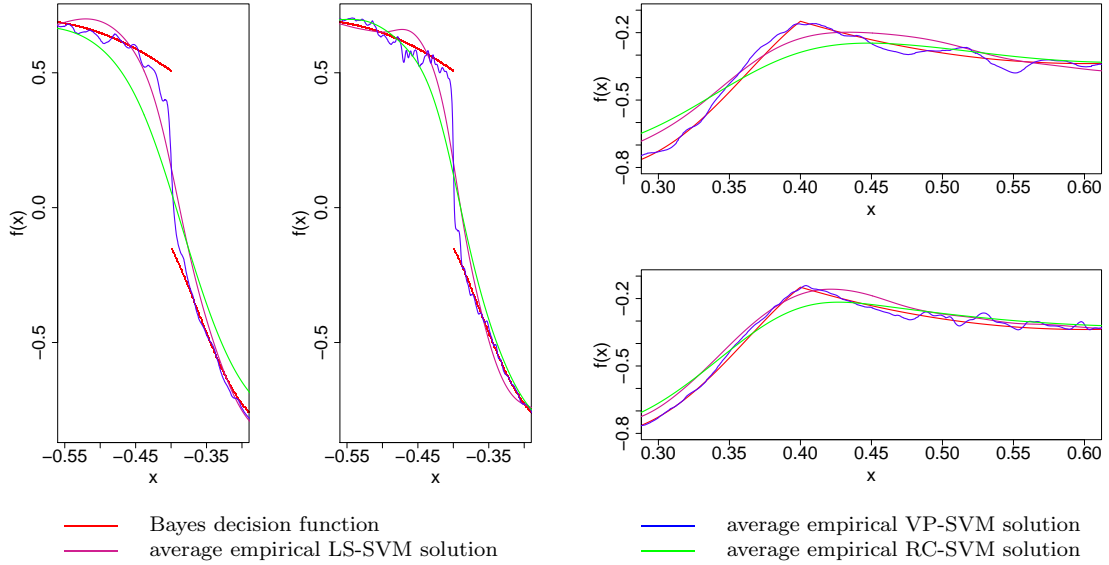


Figure 17: Predictions for the artificial data sets of Type II. The left graphic shows the predictions for $x \in [-0.55, -0.4]$ and the data set of size $n_{\text{train}} = 1000$, while the graphic on its right-hand side pictures the predictions for the same interval for x and the data set of size $n_{\text{train}} = 10000$. The two graphics on the right-hand side illustrate the predictions for $x \in [0.3, 0.6]$, the upper one for the data set of size $n_{\text{train}} = 1000$ and the lower one for the data set of size $n_{\text{train}} = 10000$.

of $f_{L,P}^*$ almost completely. Moreover, in Figure 16 the smoothness interruptions of $f_{L,P}^*$ are again best illustrated by the VP-SVMs, which becomes even more evident in Figure 17. Besides, Figure 18 illustrates that the peaks of $f_{L,P}^*$ are best reproduced by the VP-SVMs.

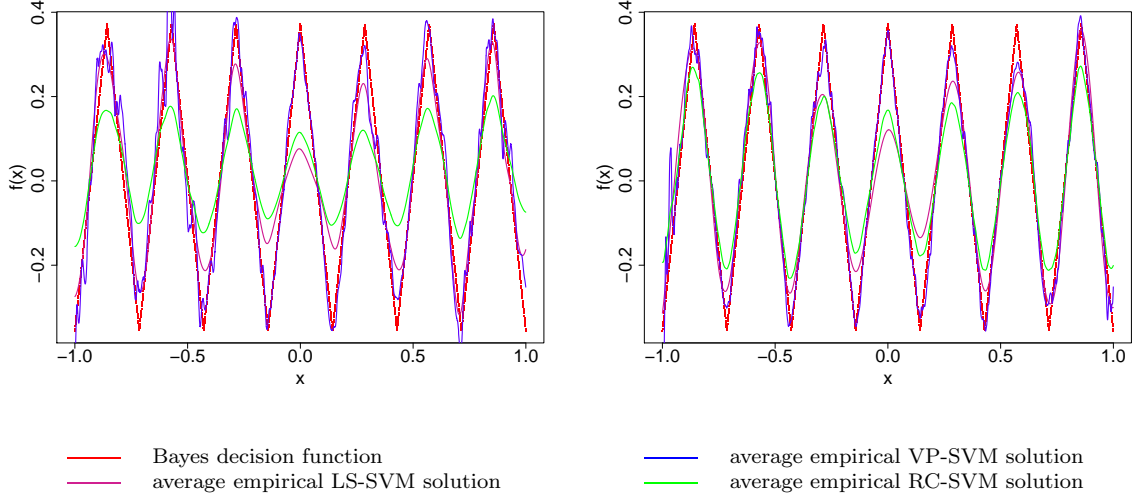


Figure 18: Predictions for the artificial data sets of Type III, drawn from the jagged function in Figure 6(c) with noise depending on x . The left graphic shows the predictions for the data set of size $n_{\text{train}} = 1000$ and the right graphic for the data set of size $n_{\text{train}} = 10000$.

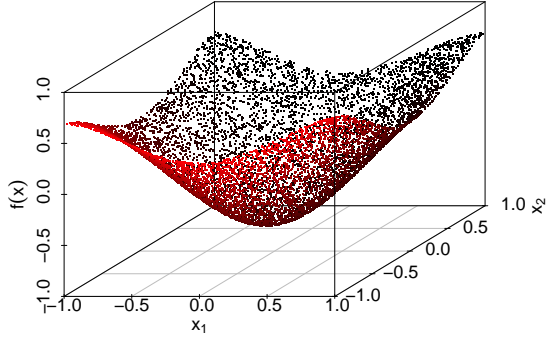
Considering the LS- and the RC-SVMs, we can not draw an universally valid conclusion, which one performs worse. In particular, Figures 15 and 16 show that, for the artificial data sets of Type I and II, both, the average empirical LS- and RC-SVM solutions, are not very well suited to the Bayes decision function. Considering the data sets of Type III, the LS-SVMs dominate the RC-SVMs in terms of the better test errors, though both kinds of SVMs do not reproduce the peaks of $f_{L,P}^*$, especially for small values of $|x|$.

It remains to optically analyze the results of the two-dimensional data sets in the following. For the artificial data sets of Type IV, the VP-SVM using 10 000 training samples achieves the best test error. Moreover, ensuing the optical impression, this VP-SVM is the only one of the considered SVM types that reflects the circular steps of the Bayes decision function as in Figure 6(d), cf. Figure 19. Finally, for the data sets of Type V, it is always the LS-SVM which performs best in terms of the test errors, cf. Table 11. This observation is also substantiated optically. To be more precise, for $n_{\text{train}} = 1000$, the uneven average empirical decision function induced by the VP-SVM (cf. Figure 20) shows that the RC-SVM even performs better than the VP-SVM. However, for $n_{\text{train}} = 10000$, the VP-SVM results are substantially improved such that the RC-SVM is now outperformed by the VP-SVM.

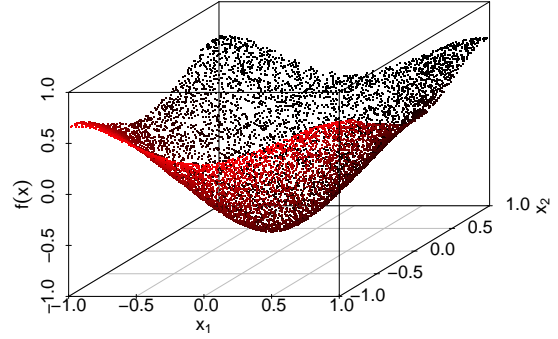
Recapulatory, we realize that the VP-SVMs possess the most distinctive ability to handle smoothness interruptions of the Bayes decision function in most of our artificial data cases, especially if $n_{\text{train}} = 10000$. For the sake of completeness, we point out that the worst performance was induced by the RC-SVMs in almost all cases, in particular for a training sample size amounting to 10 000.

6.3 Conclusions

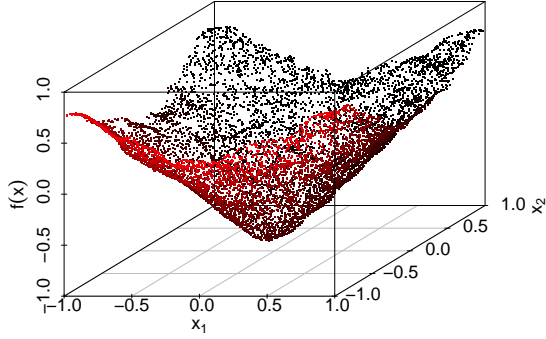
Finally, we summarize the essential findings of the previous subsections, where we considered standard LS-SVMs and two kinds of localized SVMs, namely VP-SVMs and RC-SVMs. As



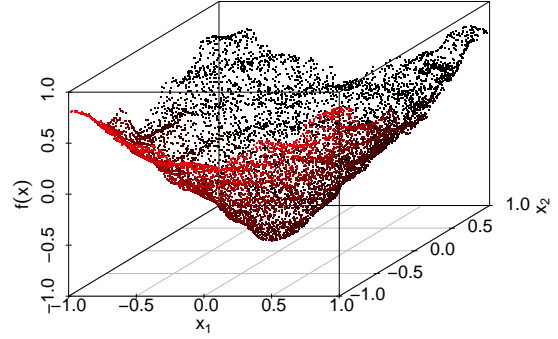
(a) LS-SVM using 1000 training samples



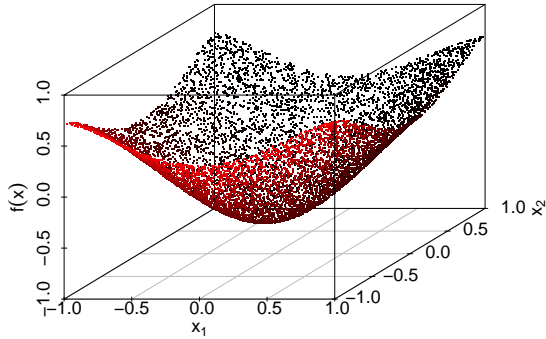
(b) LS-SVM using 10000 training samples



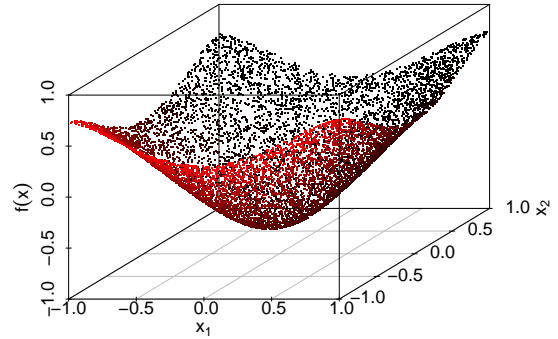
(c) VP-SVM using radius $r = 0.5$ and 1000 training samples



(d) VP-SVM using radius $r = 0.5$ and 10000 training samples

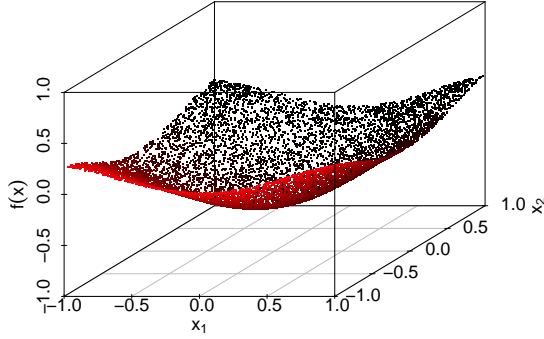


(e) RC-SVM using 20 working sets and 1000 training samples

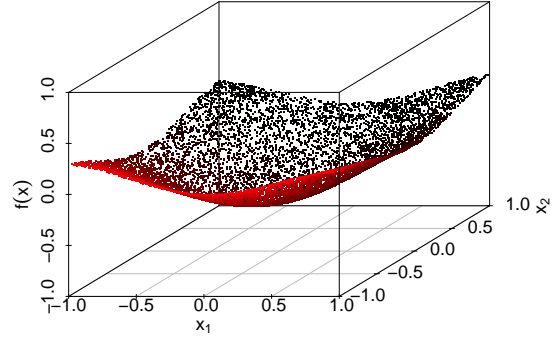


(f) RC-SVM using 20 working sets and 10000 training samples

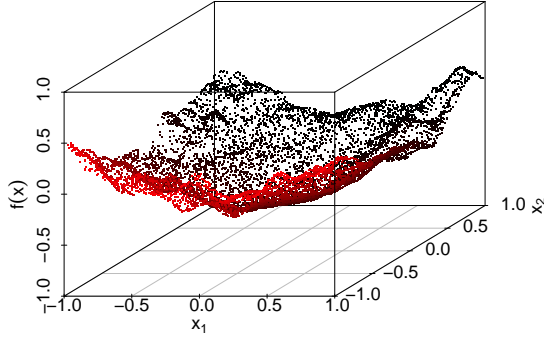
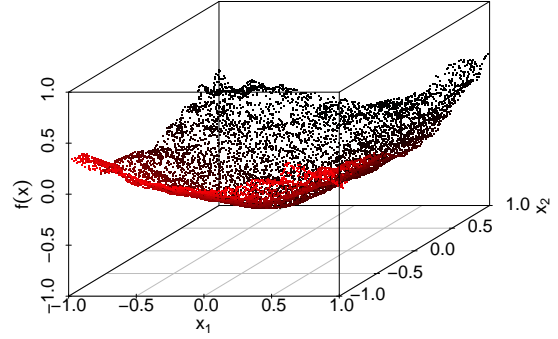
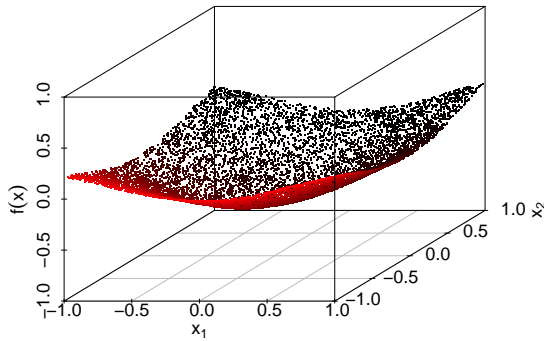
Figure 19: Predictions for the artificial data sets of Type IV, drawn from the circular step function in Figure 6(d) with noise independent of x .



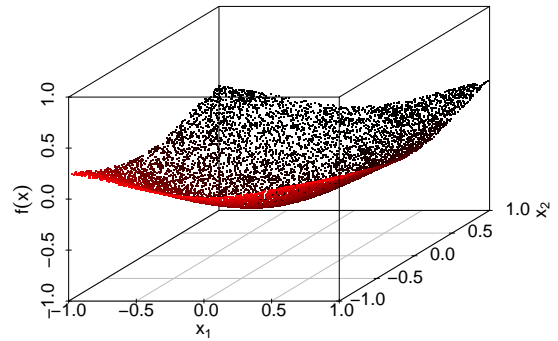
(a) LS-SVM using 1000 training samples



(b) LS-SVM using 10000 training samples


 (c) VP-SVM using radius $r = 0.5$ and 1000 training samples

 (d) VP-SVM using radius $r = 0.5$ and 10000 training samples


(e) RC-SVM using 15 working sets and 1000 training samples



(f) RC-SVM using 15 working sets and 10000 training samples

 Figure 20: Predictions for the artificial data sets of Type V, drawn from the 2-dimensional Euclidean norm in Figure 6(e) with noise independent of x .

just analyzed in Subsection 6.2, VP-SVMs have the evident advantage that they manage smoothness interruptions of the Bayes decision function better than LS- and RC-SVMs.

The real-world data sets demonstrated that the RC-SVMs perform considerably worse than the LS-SVMs and the VP-SVMs, while the performance of VP-SVMs using small radii is improved for increasing sample sizes. To be more precise, VP-SVMs outperform LS-SVMs or at most leads to a negligible worsening compared to LS-SVMs for a fraction of the training time and without memory constraints on the large data sets. For very small data sets, however, LS-SVMs actually train faster than VP-SVMs and, hence, are preferable. What is more, for data sets of size $n_{\text{train}} \leq 2500$ all LS-SVMs require less than 9s to train, so that there are probably no reasons to apply a VP-SVM. Besides, really small training sample sizes involve considerably smaller working sets for a VP-SVM using a small radius, so that it is hard to find a well suited prediction.

Furthermore, despite a faster training procedure, a VP-SVM using a sufficiently small radius induces considerably lower test errors for sample sizes $n_{\text{train}} > 100\,000$ than a LS-SVM for training data sets that still enable computational feasibility.

7. Proofs

This section is dedicated to prove the results of the previous sections. We begin with the proof of Lemma 1 relating the radius r of a cover $B_r(z_1), \dots, B_r(z_m)$ of X defined by (6) with the number m of centers z_1, \dots, z_m .

Proof [Proof of Lemma 1] First of all, let us recall the m -th entropy number of X defined by

$$\varepsilon_m(X) := \inf \left\{ \varepsilon > 0 : \exists z_1, \dots, z_m \in X \text{ such that } X \subset \bigcup_{j=1}^m (z_j + \varepsilon B_{\ell_2^d}) \right\}.$$

Since $X \subset cB_{\ell_2^d}$, the m -th entropy number of X can be upper bounded by

$$\varepsilon_m(X) \leq 2\varepsilon_m(cB_{\ell_2^d}) \leq 2c\varepsilon_m(B_{\ell_2^d}).$$

Additionally, we know by (Carl and Stephani, 1990, Section 1.1) that

$$m^{-\frac{1}{d}} \leq \varepsilon_m(B_{\ell_2^d}) \leq 4m^{-\frac{1}{d}},$$

so that we can find a cover $(B_j)_{j=1, \dots, m}$ of $X \subset cB_{\ell_2^d}$ satisfying

$$r \leq 8cm^{-\frac{1}{d}}.$$

■

7.1 Proofs of Section 3

In Section 3 we presented a lemma that related the risk w.r.t. the loss L to the risk w.r.t. the restricted loss L_j and also transferred this result to the excess risk. Hereafter, the proof of this lemma can be found.

Proof [Proof of Lemma 4] Simple transformations using $A \cup B = X$ and $A \cap B = \emptyset$ show

$$\begin{aligned}\mathcal{R}_{L,P}(f) &= \int_{X \times Y} L(x, y, \mathbb{1}_A(x)f_A(x) + \mathbb{1}_B(x)f_B(x)) dP(x, y) \\ &= \int_{X \times Y} \mathbb{1}_A(x)L(x, y, f_A(x)) + \mathbb{1}_B(x)L(x, y, f_B(x)) dP(x, y) \\ &= \mathcal{R}_{L_A,P}(f_A) + \mathcal{R}_{L_B,P}(f_B).\end{aligned}$$

The second assertion follows immediately. ■

To derive the new oracle inequality of Theorem 5 we first have to relate the entropy numbers of H_j , $j \in \{1, \dots, m\}$, to those of H . To this end, we consider a similar concept to entropy numbers, namely covering numbers, cf. (Györfi et al., 2002, Definition 9.3) or (Steinwart and Christmann, 2008a, Definition 6.19).

Definition 11 Let (T, d) be a metric space and $\varepsilon > 0$. A subset $S \subset T$ is called an ε -net of T if for all $t \in T$ there exists an $s \in S$ such that $d(s, t) \leq \varepsilon$. Furthermore, we define the ε -covering number of T by

$$\mathcal{N}(T, d, \varepsilon) := \inf \left\{ n \geq 1 : \exists s_1, \dots, s_n \in T \text{ such that } T \subset \bigcup_{i=1}^n B_d(s_i, \varepsilon) \right\},$$

where $\inf \emptyset := \infty$ and $B_d(s, \varepsilon) := \{t \in T : d(t, s) \leq \varepsilon\}$.

Note that an upper bound on entropy numbers involves a bound on covering numbers. To be more precise, for a metric space (T, d) and constants $a > 0$ and $q > 0$, the implication

$$e_i(T, d) \leq ai^{-1/q}, \quad i \geq 1 \quad \implies \quad \ln \mathcal{N}(T, d, \varepsilon) \leq \ln(4) \left(\frac{a}{\varepsilon} \right)^q, \quad \forall \varepsilon > 0 \quad (26)$$

holds by (Steinwart and Christmann, 2008a, Lemma 6.21). Additionally, (Steinwart and Christmann, 2008a, Exercise 6.8) yields the opposite implication, namely

$$\ln \mathcal{N}(T, d, \varepsilon) < \left(\frac{a}{\varepsilon} \right)^q, \quad \varepsilon > 0 \quad \implies \quad e_i(T, d) \leq 3^{1/q} ai^{-1/q}, \quad \forall i \geq 1. \quad (27)$$

Recall that we pursue the target to estimate $e_i(\text{id} : H \rightarrow L_2(P_X))$. In fact, the equivalence of entropy and covering numbers enables us to estimate the covering number $\mathcal{N}(B_H, \|\cdot\|_{L_2(P_X)}, \varepsilon)$ of H instead.

Lemma 12 Let ν be a distribution on X and $A, B \subset X$ with $A \cap B = \emptyset$. Moreover, let H_A and H_B be RKHSs on A and B that are embedded into $L_2(\nu|_A)$ and $L_2(\nu|_B)$, respectively. Let the extended RKHSs \hat{H}_A and \hat{H}_B be defined as in Lemma 2 and denote their direct sum by H as in (13), where the norm is given by (14) with $\lambda_A, \lambda_B > 0$. Then, for the ε -covering number of H w.r.t. $\|\cdot\|_{L_2(\nu)}$, we have

$$\mathcal{N}(B_H, \|\cdot\|_{L_2(\nu)}, \varepsilon) \leq \mathcal{N}\left(\lambda_A^{-1/2} B_{\hat{H}_A}, \|\cdot\|_{L_2(\nu|_A)}, \varepsilon_A\right) \cdot \mathcal{N}\left(\lambda_B^{-1/2} B_{\hat{H}_B}, \|\cdot\|_{L_2(\nu|_B)}, \varepsilon_B\right),$$

where $\varepsilon_A, \varepsilon_B > 0$ and $\varepsilon := \sqrt{\varepsilon_A^2 + \varepsilon_B^2}$.

Proof First of all, we assume that there exist $a, b \in \mathbb{N}$ and functions $\hat{f}_1, \dots, \hat{f}_a \in \lambda_A^{-\frac{1}{2}} B_{\hat{H}_A}$ and $\hat{h}_1, \dots, \hat{h}_b \in \lambda_B^{-\frac{1}{2}} B_{\hat{H}_B}$ such that $\{\hat{f}_1, \dots, \hat{f}_a\}$ is an ε_A -cover of $\lambda_A^{-\frac{1}{2}} B_{\hat{H}_A}$ w.r.t. $\|\cdot\|_{L_2(\nu|_A)}$, $\{\hat{h}_1, \dots, \hat{h}_b\}$ is an ε_B -cover of $\lambda_B^{-\frac{1}{2}} B_{\hat{H}_B}$ w.r.t. $\|\cdot\|_{L_2(\nu|_B)}$,

$$a = \mathcal{N}(\lambda_A^{-\frac{1}{2}} B_{\hat{H}_A}, \|\cdot\|_{L_2(\nu|_A)}, \varepsilon_A) \quad \text{and} \quad b = \mathcal{N}(\lambda_B^{-\frac{1}{2}} B_{\hat{H}_B}, \|\cdot\|_{L_2(\nu|_B)}, \varepsilon_B).$$

That is, for every function $\hat{g}_A \in \lambda_A^{-\frac{1}{2}} B_{\hat{H}_A}$, there exists an $i_A \in \{1, \dots, a\}$ such that

$$\|\hat{g}_A - \hat{f}_{i_A}\|_{L_2(\nu|_A)} \leq \varepsilon_A, \quad (28)$$

and for every function $\hat{g}_B \in \lambda_B^{-\frac{1}{2}} B_{\hat{H}_B}$, there exists an $i_B \in \{1, \dots, b\}$ such that

$$\|\hat{g}_B - \hat{h}_{i_B}\|_{L_2(\nu|_B)} \leq \varepsilon_B. \quad (29)$$

Let us now consider an arbitrary function $g \in B_H$. Then there exists an $\hat{g}_A \in \lambda_A^{-\frac{1}{2}} B_{\hat{H}_A}$ and an $\hat{g}_B \in \lambda_B^{-\frac{1}{2}} B_{\hat{H}_B}$ such that $g = \hat{g}_A + \hat{g}_B$. Together with (28) and (29), this implies

$$\begin{aligned} \|g - (\hat{f}_{i_A} + \hat{h}_{i_B})\|_{L_2(\nu)}^2 &= \|(\hat{g}_A - \hat{f}_{i_A}) + (\hat{g}_B - \hat{h}_{i_B})\|_{L_2(\nu)}^2 \\ &= \|\hat{g}_A - \hat{f}_{i_A}\|_{L_2(\nu|_A)}^2 + \|\hat{g}_B - \hat{h}_{i_B}\|_{L_2(\nu|_B)}^2 \\ &\leq \varepsilon_A^2 + \varepsilon_B^2 \\ &=: \varepsilon^2. \end{aligned}$$

With this, we know that

$$\left\{ \hat{f}_{i_A} + \hat{h}_{i_B} : \hat{f}_{i_A} \in \{\hat{f}_1, \dots, \hat{f}_a\} \text{ and } \hat{h}_{i_B} \in \{\hat{h}_1, \dots, \hat{h}_b\} \right\}$$

is an ε -net of H w.r.t. $\|\cdot\|_{L_2(\nu)}$. Concerning the ε -covering number of H , this finally implies

$$\mathcal{N}(B_H, \|\cdot\|_{L_2(\nu)}, \varepsilon) \leq a \cdot b = \mathcal{N}(\lambda_A^{-1/2} B_{\hat{H}_A}, \|\cdot\|_{L_2(\nu|_A)}, \varepsilon_A) \cdot \mathcal{N}(\lambda_B^{-1/2} B_{\hat{H}_B}, \|\cdot\|_{L_2(\nu|_B)}, \varepsilon_B).$$

■

Based on Lemma 12, the following theorem relates entropy numbers of H_A and H_B to those of H .

Theorem 13 *Let P_X be a distribution on X and $A_1, \dots, A_m \subset X$ be pairwise disjoint. Moreover, we assume (H) with weights $\lambda_1, \dots, \lambda_m > 0$. In addition, assume that there exist constants $p \in (0, 1)$ and $a_j > 0$, $j \in \{1, \dots, m\}$, such that for every $j \in \{1, \dots, m\}$*

$$e_i(\text{id} : H_j \rightarrow L_2(P_{X|A_j})) \leq a_j i^{-\frac{1}{2p}}, \quad i \geq 1. \quad (30)$$

Then we have

$$e_i(\text{id} : H \rightarrow L_2(\mathbb{P}_X)) \leq 2\sqrt{m} \left(3 \ln(4) \sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}} i^{-\frac{1}{2p}}, \quad i \geq 1,$$

and, for the average entropy numbers,

$$\mathbb{E}_{D_X \sim \mathbb{P}_X^n} e_i(\text{id} : H \rightarrow L_2(D_X)) \leq c_p \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}} i^{-\frac{1}{2p}}, \quad i, n \geq 1.$$

Proof [Proof of Theorem 13] First of all, note that the restriction operator $\mathcal{I} : B_{\hat{H}_j} \rightarrow B_{H_j}$ with $\mathcal{I}\hat{f} = f$ is an isometric isomorphism. Together with (Steinwart and Christmann, 2008a, (A.36)) and assumption (30), this yields

$$\begin{aligned} e_i(\lambda_j^{-\frac{1}{2}} B_{\hat{H}_j}, L_2(\mathbb{P}_{X|A_j})) &= 2\lambda_j^{-\frac{1}{2}} e_i(B_{\hat{H}_j}, L_2(\mathbb{P}_{X|A_j})) \\ &\leq 2\lambda_j^{-\frac{1}{2}} \|\mathcal{I} : B_{\hat{H}_j} \rightarrow B_{H_j}\| e_i(B_{H_j}, L_2(\mathbb{P}_{X|A_j})) \\ &\leq 2\lambda_j^{-\frac{1}{2}} a_j i^{-\frac{1}{2p}}. \end{aligned}$$

Furthermore, we know by (26) that

$$\ln \mathcal{N} \left(\lambda_j^{-\frac{1}{2}} B_{\hat{H}_j}, \|\cdot\|_{L_2(\mathbb{P}_{X|A_j})}, \varepsilon \right) \leq \ln(4) \left(2\lambda_j^{-\frac{1}{2}} a_j \right)^{2p} \varepsilon^{-2p}$$

holds for all $\varepsilon > 0$. With this and $\varepsilon_j := \frac{\varepsilon}{\sqrt{m}}$ for every $j \in \{1, \dots, m\}$, Lemma 12 implies

$$\begin{aligned} \ln \mathcal{N}(B_H, \|\cdot\|_{L_2(\mathbb{P}_X)}, \varepsilon) &\leq \ln \left(\prod_{j=1}^m \mathcal{N} \left(\lambda_j^{-\frac{1}{2}} B_{\hat{H}_j}, \|\cdot\|_{L_2(\mathbb{P}_{X|A_j})}, \varepsilon_j \right) \right) \\ &= \sum_{j=1}^m \ln \mathcal{N} \left(\lambda_j^{-\frac{1}{2}} B_{\hat{H}_j}, \|\cdot\|_{L_2(\mathbb{P}_{X|A_j})}, \frac{\varepsilon}{\sqrt{m}} \right) \\ &\leq \sum_{j=1}^m \ln(4) \left(2\lambda_j^{-\frac{1}{2}} a_j \right)^{2p} \left(\frac{\sqrt{m}}{\varepsilon} \right)^{2p} \\ &= \left(2 \ln(4)^{\frac{1}{2p}} \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}} \right)^{2p} \varepsilon^{-2p}. \end{aligned}$$

Using (27), the latter bound for the covering number of B_H finally implies the following entropy estimate

$$e_i(\text{id} : H \rightarrow L_2(\mathbb{P}_X)) \leq 3^{\frac{1}{2p}} \left(2 \ln(4)^{\frac{1}{2p}} \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}} \right) i^{-\frac{1}{2p}}$$

$$\leq 2(3\ln(4))^{\frac{1}{2p}} \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}} i^{-\frac{1}{2p}}.$$

The second assertion immediately follows by (Steinwart and Christmann, 2008a, Corollary 7.31). \blacksquare

Applying Theorem 13, we now prove Theorem 5 and thus an oracle inequality for VP-SVMs using an ordinary type of losses.

Proof [Proof of Theorem 5] Since H_1, \dots, H_m are separable RKHS of measurable kernels k_1, \dots, k_m , H is a separable RKHS and its kernel k is measurable, too. Furthermore, Theorem 13 yields

$$\mathbb{E}_{D_X \sim P_X^n} e_i(\text{id} : H \rightarrow L_2(D_X)) \leq c_p \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}} i^{-\frac{1}{2p}}, \quad i, n \geq 1.$$

That is, we can apply (Steinwart and Christmann, 2008a, Theorem 7.23) for a regularization parameter $\tilde{\lambda} = 1$ and, for all fixed $\tau > 0$ and $\lambda_j > 0$, $j \in \{1, \dots, m\}$, we obtain

$$\begin{aligned} & \sum_{j=1}^m \lambda_j \|f_{D_j, \lambda_j}\|_{\hat{H}_j}^2 + \mathcal{R}_{L_J, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L_J, P}^* \\ &= \|f_{D, \lambda}\|_H^2 + \mathcal{R}_{L_J, P}(\hat{f}_{D, \lambda}) - \mathcal{R}_{L_J, P}^* \\ &\leq 9(\|f_0\|_H^2 + \mathcal{R}_{L_J, P}(f_0) - \mathcal{R}_{L_J, P}^*) + C(a^{2p}n^{-1})^{\frac{1}{2-p-\vartheta+\vartheta p}} + 3\left(\frac{72V\tau}{n}\right)^{\frac{1}{2-\vartheta}} + \frac{15B_0\tau}{n} \\ &\leq 9\left(\sum_{j=1}^m \lambda_j \|\mathbb{1}_{A_j} f_0\|_{\hat{H}_j}^2 + \mathcal{R}_{L_J, P}(f_0) - \mathcal{R}_{L_J, P}^*\right) + C(a^{2p}n^{-1})^{\frac{1}{2-p-\vartheta+\vartheta p}} + 3\left(\frac{72V\tau}{n}\right)^{\frac{1}{2-\vartheta}} + \frac{15B_0\tau}{n} \end{aligned}$$

with probability P^n not less than $1 - 3e^{-\tau}$, where $C > 0$ is the constant of (Steinwart and Christmann, 2008a, Theorem 7.23) only depending on p , M , V , ϑ , and B . Moreover,

$$a := \max \left\{ c_p \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}}, B \right\},$$

where we need $a \geq B$ since it is a condition of (Steinwart and Christmann, 2008a, Theorem 7.23). \blacksquare

7.2 Proofs Related to the Entropy Estimates of Section 4

In this subsection, just as in Section 4, we focus on Gaussian RBF kernels and the associated RKHSs. To be more precise, we derive a bound for the entropy numbers of $H_\gamma(A)$, where $\gamma > 0$ and $A \subset \mathbb{R}^d$ with $\mathring{A} \neq \emptyset$.

Proof [Proof of Theorem 6] First of all, we consider the commutative diagram

$$\begin{array}{ccc}
 H_\gamma(A) & \xrightarrow{\text{id}} & L_2(\mathbb{P}_{X|A}) \\
 \mathcal{I}_B^{-1} \circ \mathcal{I}_A \downarrow & & \uparrow \text{id} \\
 H_\gamma(B) & \xrightarrow{\text{id}} & \ell_\infty(B)
 \end{array}$$

where the extension operator $\mathcal{I}_A : H_\gamma(A) \rightarrow H_\gamma(\mathbb{R}^d)$ and the restriction operator $\mathcal{I}_B^{-1} : H_\gamma(\mathbb{R}^d) \rightarrow H_\gamma(B)$ given by (Steinwart and Christmann, 2008a, Corollary 4.43) are isometric isomorphisms, so that $\|\mathcal{I}_B^{-1} \circ \mathcal{I}_A : H_\gamma(A) \rightarrow H_\gamma(B)\| = 1$. Furthermore, for $f \in \ell_\infty(B)$, we have

$$\|f\|_{L_2(\mathbb{P}_{X|A})} = \left(\int_X \mathbb{1}_A |f(x)|^2 d\mathbb{P}_X(x) \right)^{\frac{1}{2}} \leq \|f\|_\infty \left(\int_X \mathbb{1}_A d\mathbb{P}_X(x) \right)^{\frac{1}{2}} = \sqrt{\mathbb{P}_X(A)} \|f\|_\infty,$$

i.e. $\|\text{id} : \ell_\infty(B) \rightarrow L_2(\mathbb{P}_{X|A})\| \leq \sqrt{\mathbb{P}_X(A)}$. Together with (Steinwart and Christmann, 2008a, (A.38) and (A.39)) as well as (Steinwart and Christmann, 2008a, Theorem 6.27), we obtain for all $i \geq 1$

$$\begin{aligned}
 & e_i(\text{id} : H_\gamma(A) \rightarrow L_2(\mathbb{P}_{X|A})) \\
 & \leq \|\mathcal{I}_B^{-1} \circ \mathcal{I}_A : H_\gamma(A) \rightarrow H_\gamma(B)\| \cdot e_i(\text{id} : H_\gamma(B) \rightarrow \ell_\infty(B)) \cdot \|\text{id} : \ell_\infty(B) \rightarrow L_2(\mathbb{P}_{X|A})\| \\
 & \leq \sqrt{\mathbb{P}_X(A)} c_{m,d} r^m \gamma^{-m} i^{-\frac{m}{d}},
 \end{aligned}$$

where $m \geq 1$ is an arbitrary integer and $c_{m,d}$ a positive constant. For $p \in (0, 1)$, the choice $m = \left\lceil \frac{d}{2p} \right\rceil$ finally yields

$$e_i(\text{id} : H_\gamma(A) \rightarrow L_2(\mathbb{P}_{X|A})) \leq \sqrt{\mathbb{P}_X(A)} c_{m,d} r^m \gamma^{-m} i^{-\frac{m}{d}} \leq c_p \sqrt{\mathbb{P}_X(A)} r^{\frac{d+2p}{2p}} \gamma^{-\frac{d+2p}{2p}} i^{-\frac{1}{2p}}.$$

■

7.3 Proofs Related to the Least Squares VP-SVMs

In this subsection, we prove the results that are linked with the least squares loss, i.e. the results of Section 5. Before we elaborate on the oracle inequality for VP-SVMs using the least squares loss as well as RKHSs of Gaussian kernels, we have to examine the excess risk

$$\mathcal{R}_{L_{J_T}, \mathbb{P}}(f_0) - \mathcal{R}_{L_{J_T}, \mathbb{P}}^* = \|f_0 - f_{L, \mathbb{P}}^*\|_{L_2(\mathbb{P}_{X|A_T})}^2. \quad (31)$$

Let us begin by writing for fixed $\gamma_j > 0$

$$K_j : \mathbb{R}^d \rightarrow \mathbb{R}, \quad x \mapsto \sum_{\ell=1}^s \binom{s}{\ell} (-1)^{1-\ell} \left(\frac{2}{\ell^2 \gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|x\|_2^2}{\ell^2 \gamma_j^2} \right), \quad (32)$$

and choosing $f_0 := \sum_{j=1}^m \mathbb{1}_{A_j} \cdot (K_j * f_{L,P}^*)$. Then (31) can be estimated with the help of the following theorem, which is together with its proof basically a modification of (Eberts and Steinwart, 2013, Theorem 2.2). Indeed, the proofs proceed mainly identically. Note that we use the notation

$$\gamma_{\max} := \max\{\gamma_1, \dots, \gamma_m\} \quad \text{and} \quad \gamma_{\min} := \min\{\gamma_1, \dots, \gamma_m\}$$

in the following theorem and the associated proof.

Theorem 14 *Let us fix some $q \in [1, \infty)$. Assume that ν is a finite measure on \mathbb{R}^d with $\text{supp } \nu =: X \subset \mathbb{R}^d$ and let $(A_j)_{j=1, \dots, m}$ be a partition of X . Furthermore, let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be such that $f \in B_{q, \infty}^\alpha(\nu)$ for some $\alpha \geq 1$. For the functions $K_j : \mathbb{R}^d \rightarrow \mathbb{R}$, $j \in \{1, \dots, m\}$, defined by (32), where $s := \lfloor \alpha \rfloor + 1$ and $\gamma_1, \dots, \gamma_m > 0$, we then have*

$$\left\| \sum_{j=1}^m \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_q(\nu)}^q \leq C_{\alpha, q} \left(\frac{\gamma_{\max}}{\gamma_{\min}} \right)^d \gamma_{\max}^{q\alpha},$$

where $C_{\alpha, q} := \|f\|_{B_{q, \infty}^\alpha(\nu)}^q \left(\frac{d}{2} \right)^{\frac{q\alpha}{2}} \pi^{-\frac{1}{4}} \Gamma\left(q\alpha + \frac{1}{2}\right)^{\frac{1}{2}}$.

Proof In the following, we write $J := \{1, \dots, m\}$. To show

$$\left\| \sum_{j \in J} \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_q(\nu)}^q \leq \|f\|_{B_{q, \infty}^\alpha(\nu)}^q \left(\frac{d}{2} \right)^{\frac{q\alpha}{2}} \pi^{-\frac{1}{4}} \Gamma\left(q\alpha + \frac{1}{2}\right)^{\frac{1}{2}} \left(\frac{\gamma_{\max}}{\gamma_{\min}} \right)^d \gamma_{\max}^{q\alpha},$$

we have to proceed in a similar way as in the proof of (Eberts and Steinwart, 2013, Theorem 2.2). First of all, we use the translation invariance of the Lebesgue measure and $\exp(-\|u\|_2^2) = \exp(-\| -u \|_2^2)$ ($u \in \mathbb{R}^d$) to obtain, for $x \in X$ and $j \in J$,

$$\begin{aligned} K_j * f(x) &= \int_{\mathbb{R}^d} \sum_{\ell=1}^s \binom{s}{\ell} (-1)^{1-\ell} \frac{1}{\ell^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp\left(-\frac{2\|x-t\|_2^2}{\ell^2 \gamma_j^2}\right) f(t) dt \\ &= \int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp\left(-\frac{2\|h\|_2^2}{\gamma_j^2}\right) \left(\sum_{\ell=1}^s \binom{s}{\ell} (-1)^{1-\ell} f(x + \ell h) \right) dh. \end{aligned}$$

With this we can derive, for $q \geq 1$,

$$\begin{aligned} &\left\| \sum_{j \in J} \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_q(\nu)}^q \\ &= \int_{\mathbb{R}^d} \left| \sum_{j \in J} \mathbb{1}_{A_j}(x) (K_j * f)(x) - f(x) \right|^q d\nu(x) \\ &\leq \int_{\mathbb{R}^d} \left(\sum_{j \in J} \mathbb{1}_{A_j}(x) |K_j * f(x) - f(x)| \right)^q d\nu(x) \end{aligned}$$

$$\begin{aligned}
 &= \int_{\mathbb{R}^d} \sum_{j \in J} \mathbb{1}_{A_j}(x) |K_j * f(x) - f(x)|^q d\nu(x) \\
 &= \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) |K_j * f(x) - f(x)|^q d\nu(x) \\
 &= \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) \left| \int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) \left(\sum_{\ell=0}^s \binom{s}{\ell} (-1)^{2s+1-\ell} f(x + \ell h) \right) dh \right|^q d\nu(x) \\
 &= \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) \left| \int_{\mathbb{R}^d} (-1)^{s+1} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) \triangle_h^s(f, x) dh \right|^q d\nu(x) \\
 &\leq \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) \left(\int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) |\triangle_h^s(f, x)| dh \right)^q d\nu(x).
 \end{aligned}$$

Then Hölder's inequality and $\int_{\mathbb{R}^d} \exp \left(-2\gamma_j^{-2}\|h\|_2^2 \right) dh = \left(\frac{\gamma_j^2 \pi}{2} \right)^{d/2}$ yield, for $q > 1$,

$$\begin{aligned}
 &\left\| \sum_{j \in J} \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_q(\nu)}^q \\
 &\leq \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) \left(\left(\int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) dh \right)^{\frac{q-1}{q}} \right. \\
 &\quad \left. \left(\int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) |\triangle_h^s(f, x)|^q dh \right)^{\frac{1}{q}} \right)^q d\nu(x) \\
 &= \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) \int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) |\triangle_h^s(f, x)|^q dh d\nu(x) \\
 &= \sum_{j \in J} \int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) |\triangle_h^s(f, x)|^q d\nu(x) dh \\
 &\leq \int_{\mathbb{R}^d} \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \int_{\mathbb{R}^d} \sum_{j \in J} \mathbb{1}_{A_j}(x) |\triangle_h^s(f, x)|^q d\nu(x) dh \\
 &= \int_{\mathbb{R}^d} \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \|\triangle_h^s(f, \cdot)\|_{L_q(\nu)}^q dh \\
 &\leq \int_{\mathbb{R}^d} \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \omega_{s, L_q(\nu)}^q(f, \|h\|_2) dh.
 \end{aligned}$$

Moreover, for $q = 1$, we have

$$\begin{aligned}
& \left\| \sum_{j \in J} \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_1(\nu)} \\
& \leq \sum_{j \in J} \int_{\mathbb{R}^d} \mathbb{1}_{A_j}(x) \int_{\mathbb{R}^d} \left(\frac{2}{\gamma_j^2 \pi} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_j^2} \right) |\triangle_h^s(f, x)| \, dh \, d\nu(x) \\
& \leq \int_{\mathbb{R}^d} \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \int_{\mathbb{R}^d} \sum_{j \in J} \mathbb{1}_{A_j}(x) |\triangle_h^s(f, x)| \, d\nu(x) \, dh \\
& \leq \int_{\mathbb{R}^d} \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \omega_{s, L_1(\nu)}(f, \|h\|_2) \, dh.
\end{aligned}$$

Consequently, we can proceed in the same way for all $q \geq 1$. To this end, note that the assumption $f \in B_{q, \infty}^\alpha(\nu)$ implies $\omega_{s, L_q(\nu)}(f, t) \leq \|f\|_{B_{q, \infty}^\alpha(\nu)} t^\alpha$ for $t > 0$. The latter together with Hölder's inequality yields

$$\begin{aligned}
& \left\| \sum_{j \in J} \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_q(\nu)}^q \\
& \leq \int_{\mathbb{R}^d} \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \omega_{s, L_q(\nu)}^q(f, \|h\|_2) \, dh \\
& \leq \|f\|_{B_{q, \infty}^\alpha(\nu)}^q \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \int_{\mathbb{R}^d} \|h\|_2^{q\alpha} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \, dh \\
& \leq \|f\|_{B_{q, \infty}^\alpha(\nu)}^q \left(\frac{2}{\pi \gamma_{\min}^2} \right)^{\frac{d}{2}} \left(\int_{\mathbb{R}^d} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \, dh \right)^{\frac{1}{2}} \left(\int_{\mathbb{R}^d} \|h\|_2^{2q\alpha} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \, dh \right)^{\frac{1}{2}} \\
& = \|f\|_{B_{q, \infty}^\alpha(\nu)}^q \left(\frac{2\gamma_{\max}^2}{\pi \gamma_{\min}^4} \right)^{\frac{d}{4}} \left(\int_{\mathbb{R}^d} \|h\|_2^{2q\alpha} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \, dh \right)^{\frac{1}{2}}.
\end{aligned}$$

Using the embedding constant $d^{\frac{q\alpha-1}{2q\alpha}}$ of $\ell_{2q\alpha}^d$ to ℓ_2^d , we obtain

$$\begin{aligned}
\int_{\mathbb{R}^d} \|h\|_2^{2q\alpha} \exp \left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2} \right) \, dh & \leq d^{q\alpha-1} \sum_{\ell=1}^d \int_{\mathbb{R}^d} h_\ell^{2q\alpha} \prod_{l=1}^d \exp \left(-\frac{2h_l^2}{\gamma_{\max}^2} \right) \, d(h_1, \dots, h_d) \\
& = d^{q\alpha-1} \sum_{\ell=1}^d \left(\frac{\gamma_{\max}^2 \pi}{2} \right)^{\frac{d-1}{2}} \int_{\mathbb{R}} h_\ell^{2q\alpha} \exp \left(-\frac{2h_\ell^2}{\gamma_{\max}^2} \right) \, dh_\ell \\
& = 2d^{q\alpha} \left(\frac{\gamma_{\max}^2 \pi}{2} \right)^{\frac{d-1}{2}} \int_0^\infty t^{2q\alpha} \exp \left(-\frac{2t^2}{\gamma_{\max}^2} \right) \, dt.
\end{aligned}$$

for $\gamma > 0$. With the substitution $t = (\frac{1}{2}\gamma_{\max}^2 u)^{\frac{1}{2}}$, the functional equation $\Gamma(t+1) = t\Gamma(t)$ of the Gamma function Γ , and $\Gamma(\frac{1}{2}) = \sqrt{\pi}$ we further have

$$\begin{aligned} \int_0^\infty t^{2q\alpha} \exp\left(-\frac{2t^2}{\gamma_{\max}^2}\right) dt &= \frac{1}{2} \frac{\gamma_{\max}}{\sqrt{2}} \left(\frac{\gamma_{\max}^2}{2}\right)^{q\alpha} \int_0^\infty u^{(q\alpha+\frac{1}{2})-1} \exp(-u) du \\ &= \frac{1}{2} \frac{\gamma_{\max}}{\sqrt{2}} \left(\frac{\gamma_{\max}^2}{2}\right)^{q\alpha} \Gamma\left(q\alpha + \frac{1}{2}\right). \end{aligned}$$

Altogether, we finally obtain

$$\begin{aligned} &\left\| \sum_{j \in J} \mathbb{1}_{A_j} \cdot (K_j * f) - f \right\|_{L_q(\nu)}^q \\ &\leq \|f\|_{B_{q,\infty}^\alpha(\nu)}^q \left(\frac{2\gamma_{\max}^2}{\pi\gamma_{\min}^4}\right)^{\frac{d}{4}} \left(\int_{\mathbb{R}^d} \|h\|_2^{2q\alpha} \exp\left(-\frac{2\|h\|_2^2}{\gamma_{\max}^2}\right) dh\right)^{\frac{1}{2}} \\ &\leq \|f\|_{B_{q,\infty}^\alpha(\nu)}^q \left(\frac{2\gamma_{\max}^2}{\pi\gamma_{\min}^4}\right)^{\frac{d}{4}} \left(\left(\frac{d}{2}\right)^{q\alpha} \left(\frac{\pi^{d-1}}{2^d}\right)^{\frac{1}{2}} \gamma_{\max}^{2q\alpha+d} \Gamma\left(q\alpha + \frac{1}{2}\right)\right)^{\frac{1}{2}} \\ &= \|f\|_{B_{q,\infty}^\alpha(\nu)}^q \left(\frac{d}{2}\right)^{\frac{q\alpha}{2}} \pi^{-\frac{1}{4}} \Gamma\left(q\alpha + \frac{1}{2}\right)^{\frac{1}{2}} \left(\frac{\gamma_{\max}}{\gamma_{\min}}\right)^d \gamma_{\max}^{q\alpha}. \end{aligned}$$

■

Based on Theorems 5 and 14, we can now show Theorem 7.

Proof [Proof of Theorem 7] First, we have to choose a function $f_0 \in H$. To this end, we define functions $K_j : \mathbb{R}^d \rightarrow \mathbb{R}$, $j \in \{1, \dots, m\}$, by (32), where $s := \lfloor \alpha \rfloor + 1$ and $\gamma_j > 0$. Then we define f_0 by convolving each K_j with the Bayes decision function $f_{L,P}^*$, that is

$$f_0(x) := \sum_{j \in J_T} \mathbb{1}_{A_j}(x) \cdot (K_j * f_{L,P}^*)(x), \quad x \in \mathbb{R}^d.$$

Now, to show that f_0 is indeed a suitable function to bound the approximation error, we first need to ensure that f_0 is contained in H . In addition, we need to derive bounds for both, the regularization term and the excess risk of f_0 . To this end, we apply (Eberts and Steinwart, 2013, Theorem 2.3) and obtain, for every $j \in J_T$,

$$(K_j * f_{L,P}^*)|_{A_j} \in H_{\gamma_j}(A_j)$$

with

$$\begin{aligned} \|\mathbb{1}_{A_j} f_0\|_{\hat{H}_{\gamma_j}(A_j)} &= \|\mathbb{1}_{A_j} (K_j * f_{L,P}^*)\|_{\hat{H}_{\gamma_j}(A_j)} \\ &= \|(K_j * f_{L,P}^*)|_{A_j}\|_{H_{\gamma_j}(A_j)} \\ &\leq (\gamma_j \sqrt{\pi})^{-\frac{d}{2}} (2^s - 1) \|f_{L,P}^*\|_{L_2(\mathbb{R}^d)}. \end{aligned}$$

This implies

$$f_0 = \sum_{j \in J_T} \underbrace{\mathbb{1}_{A_j}(K_j * f_{L,P}^*)}_{\in \hat{H}_{\gamma_j}(A_j)} \in H_{J_T}.$$

Besides, note that $0 \in \hat{H}_{\gamma_j}(A_j)$ for every $j \in \{1, \dots, m\}$ such that f_0 can be written as $f_0 = \sum_{j=1}^m f_j$, where

$$f_j := \begin{cases} \mathbb{1}_{A_j}(K_j * f_{L,P}^*), & j \in J_T, \\ 0, & j \notin J_T. \end{cases}$$

Obviously, the latter implies $f_0 \in H$. Furthermore, for $A_T := \bigcup_{j \in J_T} A_j$, (31) and Theorem 14 yield

$$\begin{aligned} \mathcal{R}_{L_{J_T}, P}(f_0) - \mathcal{R}_{L_{J_T}, P}^* &= \|f_0 - f_{L,P}^*\|_{L_2(P_{X|A_T})}^2 \\ &= \left\| \sum_{j \in J_T} \mathbb{1}_{A_j}(K_j * f_{L,P}^*) - f_{L,P}^* \right\|_{L_2(P_{X|A_T})}^2 \\ &\leq C_{\alpha,2} \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha}, \end{aligned}$$

where $C_{\alpha,2}$ is a constant only depending on α , d , and $\|f_{L,P}^*\|_{B_{2,\infty}^\alpha(P_{X|A_T})}$. To utilize Theorem 5, it remains to examine the constants B , V , ϑ , and B_0 . Since we consider the least squares loss, which can be clipped at M with $Y = [-M, M]$, the supremum bound (18) holds for $B = 4M^2$ and the variance bound (19) for $V = 16M^2$ and $\vartheta = 1$ (cf. Steinwart and Christmann, 2008a, Example 7.3). Next, we derive a bound for $\|L \circ f_0\|_\infty$ using (Eberts and Steinwart, 2013, Theorem 2.3) which provides, for every $x \in X$, the supremum bound

$$|f_0(x)| = \left| \sum_{j \in J_T} \mathbb{1}_{A_j}(x) \cdot (K_j * f_{L,P}^*)(x) \right| \leq \sum_{j \in J_T} \mathbb{1}_{A_j}(x) |K_j * f_{L,P}^*(x)| \leq (2^s - 1) \|f_{L,P}^*\|_{L_\infty(\mathbb{R}^d)}. \quad (33)$$

The latter implies

$$\begin{aligned} \|L_{J_T} \circ f_0\|_\infty &= \sup_{(x,y) \in X \times Y} |L(y, f_0(x))| \\ &\leq \sup_{(x,y) \in X \times Y} (M^2 + 2M|f_0(x)| + |f_0(x)|^2) \\ &\leq 4^s \max \left\{ M^2, \|f_{L,P}^*\|_{L_\infty(\mathbb{R}^d)}^2 \right\}, \end{aligned}$$

i.e. $B_0 := 4^s \max \{ M^2, \|f_{L,P}^*\|_{L_\infty(\mathbb{R}^d)}^2 \}$. Moreover, since Theorem 6 provides $e_i(\text{id} : H_{\gamma_j}(A_j) \rightarrow L_2(P_{X|A_j})) \leq a_j i^{-\frac{1}{2p}}$ for $i \geq 1$ with $a_j = \tilde{c}_p \sqrt{P_X(A_j)} r^{\frac{d+2p}{2p}} \gamma_j^{-\frac{d+2p}{2p}}$, we have

$$\left(\max \left\{ c_p \sqrt{m} \left(\sum_{j=1}^m \lambda_j^{-p} a_j^{2p} \right)^{\frac{1}{2p}}, B \right\} \right)^{2p}$$

$$\begin{aligned}
 &= \left(\max \left\{ c_p \tilde{c}_p \sqrt{mr}^{\frac{d+2p}{2p}} \left(\sum_{j=1}^m \left(\lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^p \right)^{\frac{1}{2p}}, B \right\} \right)^{2p} \\
 &\leq \left(\max \left\{ c_p \tilde{c}_p m^{\frac{1}{2p}} r^{\frac{d+2p}{2p}} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^{\frac{1}{2}}, B \right\} \right)^{2p} \\
 &\leq \left(\max \left\{ c_p \tilde{c}_p 8^{\frac{d}{2p}} r \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^{\frac{1}{2}}, B \right\} \right)^{2p} \\
 &\leq C_p r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^p + B^{2p} \\
 &=: a^{2p},
 \end{aligned}$$

where we used the concavity of the function $t \mapsto t^p$ for $t \geq 0$, $mr^d \leq 8^d$ by (7), and $C_p := c_p^{2p} \tilde{c}_p^{2p} 8^d$. Finally, applying Theorem 5 yields

$$\begin{aligned}
 &\mathcal{R}_{L_{J_T}, \mathbb{P}}(\hat{f}_{\mathbf{D}, \boldsymbol{\lambda}, \boldsymbol{\gamma}}) - \mathcal{R}_{L_{J_T}, \mathbb{P}}^* \\
 &\leq \sum_{j=1}^m \lambda_j \|f_{\mathbf{D}_j, \lambda_j, \gamma_j}\|_{\dot{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_{J_T}, \mathbb{P}}(\hat{f}_{\mathbf{D}, \boldsymbol{\lambda}, \boldsymbol{\gamma}}) - \mathcal{R}_{L_{J_T}, \mathbb{P}}^* \\
 &\leq 9 \left(\sum_{j=1}^m \lambda_j \|\mathbb{1}_{A_j} f_0\|_{\dot{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_{J_T}, \mathbb{P}}(f_0) - \mathcal{R}_{L_{J_T}, \mathbb{P}}^* \right) \\
 &\quad + C (a^{2p} n^{-1})^{\frac{1}{2-p-\vartheta+\vartheta p}} + 3 \left(\frac{72V\tau}{n} \right)^{\frac{1}{2-\vartheta}} + \frac{15B_0\tau}{n} \\
 &\leq 9 \left(\sum_{j \in J_T} \lambda_j (\gamma_j \sqrt{\pi})^{-d} (2^s - 1)^2 \|f_{L, \mathbb{P}}^*\|_{L_2(\mathbb{R}^d)}^2 + C_{\alpha, 2} \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} \right) \\
 &\quad + C C_p r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^p n^{-1} + C B^{2p} n^{-1} + \frac{3456 M^2 \tau}{n} \\
 &\quad + 15 \cdot 4^s \max\{M^2, \|f_{L, \mathbb{P}}^*\|_{L_\infty(\mathbb{R}^d)}^2\} \frac{\tau}{n} \\
 &\leq 9(2^s - 1)^2 \pi^{-\frac{d}{2}} \|f_{L, \mathbb{P}}^*\|_{L_2(\mathbb{R}^d)}^2 \sum_{j \in J_T} \lambda_j \gamma_j^{-d} + 9C_{\alpha, 2} \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} \\
 &\quad + C C_p r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^p n^{-1} + 16^p C M^{4p} n^{-1} \\
 &\quad + \left(3456 M^2 + 15 \cdot 4^s \max\{M^2, \|f_{L, \mathbb{P}}^*\|_{L_\infty(\mathbb{R}^d)}^2\} \right) \frac{\tau}{n}
 \end{aligned}$$

with probability \mathbf{P}^n not less than $1 - 3e^{-\tau}$. Finally, for $\hat{\tau} \geq 1$, a variable transformation implies

$$\begin{aligned} & \sum_{j=1}^m \lambda_j \|f_{D_j, \lambda_j, \gamma_j}\|_{\hat{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_{J_T}, \mathbf{P}}(\hat{f}_{D, \lambda, \gamma}) - \mathcal{R}_{L_{J_T}, \mathbf{P}}^* \\ & \leq C_{M, \alpha, p} \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} + r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbf{P}_X(A_j) \right)^p n^{-1} + \hat{\tau} n^{-1} \right) \end{aligned}$$

with probability \mathbf{P}^n not less than $1 - e^{-\hat{\tau}}$, where the constant $C_{M, \alpha, p}$ is defined by

$$\begin{aligned} C_{M, \alpha, p} := \max & \left\{ 9(2^s - 1)^2 \pi^{-\frac{d}{2}} \|f_{L, \mathbf{P}}^*\|_{L_2(\mathbb{R}^d)}^2, 9 \|f_{L, \mathbf{P}}^*\|_{B_{2, \infty}^\alpha(\mathbf{P}_{X|A_T})}^2 \left(\frac{d}{2} \right)^\alpha \pi^{-\frac{1}{4}} \Gamma \left(2\alpha + \frac{1}{2} \right)^{\frac{1}{2}}, \right. \\ & \left. 8^d C_p^{2p} \tilde{c}_p^{2p}, 16^p C M^{4p} + \left(3456 M^2 + 15 \cdot 4^s \max\{M^2, \|f_{L, \mathbf{P}}^*\|_{L_\infty(\mathbb{R}^d)}^2\} \right) (1 + \ln(3)) \right\}. \end{aligned}$$

■

Next, using the just proven oracle inequality presented in Theorem 7, we show the learning rates of Theorem 8 in only a few steps.

Proof [Proof of Theorem 8] First of all, we define sequences $\tilde{\lambda}_n := c_2 n^{-1}$ and $\tilde{\gamma}_n := c_3 n^{-\frac{1}{2\alpha+d}}$ to simplify the presentation. Then Theorem 7, $\sum_{j=1}^{m_n} \mathbf{P}_X(A_j) = 1$, and $|J_T| \leq m_n \leq 8^d r_n^{-d}$ together with $\lambda_{n,j} = r_n^d \tilde{\lambda}_n$ and $\gamma_{n,j} = \tilde{\gamma}_n$ for all $j \in \{1, \dots, m_n\}$ yield

$$\begin{aligned} & \mathcal{R}_{L_{J_T}, \mathbf{P}}(\hat{f}_{D, \lambda_n, \gamma_n}) - \mathcal{R}_{L_{J_T}, \mathbf{P}}^* \\ & \leq C_{M, \alpha, p} \left(\sum_{j \in J_T} \lambda_{n,j} \gamma_{n,j}^{-d} + \left(\frac{\max_{j \in J_T} \gamma_{n,j}}{\min_{j \in J_T} \gamma_{n,j}} \right)^d \max_{j \in J_T} \gamma_{n,j}^{2\alpha} + r_n^{2p} \left(\sum_{j=1}^{m_n} \lambda_{n,j}^{-1} \gamma_{n,j}^{-\frac{d+2p}{p}} \mathbf{P}_X(A_j) \right)^p n^{-1} + \frac{\tau}{n} \right) \\ & = C_{M, \alpha, p} \left(|J_T| r_n^d \tilde{\lambda}_n \tilde{\gamma}_n^{-d} + \tilde{\gamma}_n^{2\alpha} + r_n^{(2-d)p} \tilde{\lambda}_n^{-p} \tilde{\gamma}_n^{-(d+2p)} \left(\sum_{j=1}^{m_n} \mathbf{P}_X(A_j) \right)^p n^{-1} + \tau n^{-1} \right) \\ & \leq 8^d C_{M, \alpha, p} \left(\tilde{\lambda}_n \tilde{\gamma}_n^{-d} + \tilde{\gamma}_n^{2\alpha} + \tilde{\lambda}_n^{-p} \tilde{\gamma}_n^{-(d+2p)} r_n^{(2-d)p} n^{-1} + \tau n^{-1} \right). \end{aligned}$$

Using the choices $\tilde{\lambda}_n = c_2 n^{-1}$, $\tilde{\gamma}_n = c_3 n^{-\frac{1}{2\alpha+d}}$, as well as $r_n = c_1 n^{-\frac{1}{\beta d}}$ finally implies

$$\begin{aligned} & \mathcal{R}_{L_{J_T}, \mathbf{P}}(\hat{f}_{D, \lambda_n, \gamma_n}) - \mathcal{R}_{L_{J_T}, \mathbf{P}}^* \\ & \leq 8^d C_{M, \alpha, p} \left(\tilde{\lambda}_n \tilde{\gamma}_n^{-d} + \tilde{\gamma}_n^{2\alpha} + \tilde{\lambda}_n^{-p} \tilde{\gamma}_n^{-(d+2p)} r_n^{(2-d)p} n^{-1} + \tau n^{-1} \right) \\ & \leq \hat{C}_{M, \alpha, p} \left(n^{-1} n^{\frac{d}{2\alpha+d}} + n^{-\frac{2\alpha}{2\alpha+d}} + n^p n^{\frac{d+2p}{2\alpha+d}} n^{-\frac{(2-d)p}{\beta d}} n^{-1} + \tau n^{-1} \right) \\ & = \hat{C}_{M, \alpha, p} \left(n^{-\frac{2\alpha}{2\alpha+d}} + n^{-\frac{2\alpha}{2\alpha+d}} + n^{-\frac{2\alpha}{2\alpha+d}} \left(1 + \frac{2}{2\alpha+d} + \frac{1}{\beta} - \frac{2}{\beta d} \right) p + \tau n^{-1} \right) \end{aligned}$$

$$\leq C\tau n^{-\frac{2\alpha}{2\alpha+d}+\xi}$$

with probability P^n not less than $1 - e^{-\tau}$, where $C > 0$ is a constant and $\xi \geq \left(1 + \frac{2}{2\alpha+d} + \frac{1}{\beta} - \frac{2}{\beta d}\right)p > 0$. \blacksquare

Proof [Proof of Corollary 9] For simplicity of notation, we write λ , λ_j , γ , and γ_j instead of λ_n , $\lambda_{n,j}$, γ_n , and $\gamma_{n,j}$. Since $\bigcup_{j \in J_T} A_j \subset T^{+\delta}$ for all $n \geq n_\delta$, the assumption $f_{L,P}^* \in B_{2,\infty}^\alpha(P_{X|T^{+\delta}})$ implies

$$f_{L,P}^* \in B_{2,\infty}^\alpha(P_{X|\bigcup_{j \in J_T} A_j}).$$

With this, Theorems 7 and 8 immediately yield

$$\begin{aligned} & \mathcal{R}_{L_T,P}(\widehat{f}_{D,\lambda,\gamma}) - \mathcal{R}_{L_T,P}^* \\ & \leq \sum_{j=1}^m \lambda_j \|f_{D_j,\lambda_j,\gamma_j}\|_{\widehat{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_T,P}(\widehat{f}_{D,\lambda,\gamma}) - \mathcal{R}_{L_T,P}^* \\ & \leq \sum_{j=1}^m \lambda_j \|f_{D_j,\lambda_j,\gamma_j}\|_{\widehat{H}_{\gamma_j}(A_j)}^2 + \mathcal{R}_{L_{J_T},P}(\widehat{f}_{D,\lambda,\gamma}) - \mathcal{R}_{L_{J_T},P}^* \\ & \leq C_{M,\alpha,p} \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} + r^{2p} \left(\sum_{j=1}^m \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p n^{-1} + \frac{\tau}{n} \right) \\ & \leq C\tau n^{-\frac{2\alpha}{2\alpha+d}+\xi} \end{aligned}$$

with probability P^n not less than $1 - e^{-\tau}$, where $\xi \geq \left(1 + \frac{2}{2\alpha+d} + \frac{1}{\beta} - \frac{2}{\beta d}\right)p > 0$. Moreover, the constants $C_{M,\alpha,p} > 0$ and $C > 0$ coincide with those of Theorems 7 and 8. \blacksquare

It remains to prove Theorem 10. However, we previously have to consider the following technical lemma.

Lemma 15 *Let $d \geq 1$ and $r_n := cn^{-\frac{1}{\beta d}}$ with $\beta > 1$ and a constant $c > 0$. We fix finite subsets $\Lambda_n \subset (0, r_n^d]$ and $\Gamma_n \subset (0, r_n]$ such that Λ_n is an $(r_n^d \varepsilon_n)$ -net of $(0, r_n^d]$ and Γ_n is an δ_n -net of $(0, r_n]$ with $0 < \varepsilon_n \leq n^{-1}$, $\delta_n > 0$, $r_n^d \in \Lambda_n$, and $r_n \in \Gamma_n$. Moreover, let $J \subset \{1, \dots, m_n\}$ be an arbitrary non-empty index set and $|J| \leq m_n \leq 8^d r_n^{-d}$. Then, for all $0 < \alpha < \frac{\beta-1}{2}d$, $n \geq 1$, and all $p \in (0, 1)$ with $p \leq \frac{\beta d - 2\alpha - d}{2\alpha + d + 2}$, we have*

$$\begin{aligned} & \inf_{(\lambda_j, \gamma_j)_{j=1}^{m_n} \in (\Lambda_n \times \Gamma_n)^{m_n}} \left(\sum_{j \in J} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J} \gamma_j}{\min_{j \in J} \gamma_j} \right)^d \max_{j \in J} \gamma_j^{2\alpha} + r_n^{2p} \left(\sum_{j=1}^{m_n} \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p n^{-1} \right) \\ & \leq C \left(n^{-\frac{2\alpha}{2\alpha+d}+\xi} + \delta_n^{2\alpha} \right), \end{aligned}$$

where $\xi := \left(\frac{2\alpha(2\alpha+d+2)}{(2\alpha+d)((2\alpha+d)(1+p)+2p)} + \max\left\{\frac{d-2}{\beta d}, 0\right\} \right)p$ and $C > 0$ is a constant independent of n , Λ_n , ε_n , Γ_n , and δ_n .

Proof Without loss of generality, we may assume that Λ_n and Γ_n are of the form $\Lambda_n = \{\lambda^{(1)}, \dots, \lambda^{(u)}\}$ and $\Gamma_n = \{\gamma^{(1)}, \dots, \gamma^{(v)}\}$ with $\lambda^{(u)} = r_n^d$ and $\gamma^{(v)} = r_n$ as well as $\lambda^{(i-1)} < \lambda^{(i)}$ and $\gamma^{(\ell-1)} < \gamma^{(\ell)}$ for all $i = 2, \dots, u$ and $\ell = 2, \dots, v$. With $\lambda^{(0)} := 0$ and $\gamma^{(0)} := 0$ it is easy to see that

$$\lambda^{(i)} - \lambda^{(i-1)} \leq 2r_n^d \varepsilon_n \quad \text{and} \quad \gamma^{(\ell)} - \gamma^{(\ell-1)} \leq 2\delta_n \quad (34)$$

hold for all $i = 1, \dots, u$ and $\ell = 1, \dots, v$. Furthermore, define $\lambda^* := n^{-\frac{2\alpha+d}{(2\alpha+d)(1+p)+2p}}$ and $\gamma^* := cn^{-\frac{1}{(2\alpha+d)(1+p)+2p}}$. Then there exist indices $i \in \{1, \dots, u\}$ and $\ell \in \{1, \dots, v\}$ with $\lambda^{(i-1)} \leq r_n^d \lambda^* \leq \lambda^{(i)}$ and $\gamma^{(\ell-1)} \leq \gamma^* \leq \gamma^{(\ell)}$. Together with (34), this yields

$$r_n^d \lambda^* \leq \lambda^{(i)} \leq r_n^d \lambda^* + 2r_n^d \varepsilon_n \quad \text{and} \quad \gamma^* \leq \gamma^{(\ell)} \leq \gamma^* + 2\delta_n. \quad (35)$$

Moreover, the definition of λ^* implies $\varepsilon_n \leq \lambda^*$ and the one of γ^* implies $\gamma^* \leq r_n$ for $\alpha < \frac{\beta-1}{2}d$ and $p \in (0, p^*]$, where $p^* := \frac{\beta d - 2\alpha - d}{2\alpha + d + 2}$. Additionally, it is easy to check that

$$\lambda^* (\gamma^*)^{-d} + (\gamma^*)^{2\alpha} + (\lambda^*)^{-p} (\gamma^*)^{-(d+2p)} r_n^{(2-d)p} n^{-1} \leq \hat{c} n^{-\frac{2\alpha}{(2\alpha+d)(1+p)+2p} + \max\{\frac{d-2}{\beta d}, 0\}p}, \quad (36)$$

where \hat{c} is a positive constant. Using (35), the bound $|J| \leq m_n \leq 8^d r_n^{-d}$, and (36), we obtain

$$\begin{aligned} & \inf_{(\lambda_j, \gamma_j)_{j=1}^{m_n} \in (\Lambda_n \times \Gamma_n)^{m_n}} \left(\sum_{j \in J} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J} \gamma_j}{\min_{j \in J} \gamma_j} \right)^d \max_{j \in J} \gamma_j^{2\alpha} + r_n^{2p} \left(\sum_{j=1}^{m_n} \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^p n^{-1} \right) \\ & \leq \sum_{j \in J} \lambda^{(i)} (\gamma^{(\ell)})^{-d} + (\gamma^{(\ell)})^{2\alpha} + \left(\sum_{j=1}^{m_n} (\lambda^{(i)})^{-1} (\gamma^{(\ell)})^{-\frac{d+2p}{p}} \mathbb{P}_X(A_j) \right)^p r_n^{2p} n^{-1} \\ & \leq |J| \lambda^{(i)} (\gamma^{(\ell)})^{-d} + (\gamma^{(\ell)})^{2\alpha} + (\lambda^{(i)})^{-p} (\gamma^{(\ell)})^{-(d+2p)} r_n^{2p} n^{-1} \\ & \leq |J| (r_n^d \lambda^* + 2r_n^d \varepsilon_n) (\gamma^*)^{-d} + (\gamma^* + 2\delta_n)^{2\alpha} + (r_n^d \lambda^*)^{-p} (\gamma^*)^{-(d+2p)} r_n^{2p} n^{-1} \\ & \leq 8^d \cdot 3\lambda^* (\gamma^*)^{-d} + (\gamma^* + 2\delta_n)^{2\alpha} + (\lambda^*)^{-p} (\gamma^*)^{-(d+2p)} r_n^{(2-d)p} n^{-1} \\ & \leq \tilde{c} \left(\lambda^* (\gamma^*)^{-d} + (\gamma^*)^{2\alpha} + (\lambda^*)^{-p} (\gamma^*)^{-(d+2p)} r_n^{(2-d)p} n^{-1} \right) + \tilde{c} \delta_n^{2\alpha} \\ & \leq \tilde{c} \hat{c} n^{-\frac{2\alpha}{(2\alpha+d)(1+p)+2p} + \max\{\frac{d-2}{\beta d}, 0\}p} + \tilde{c} \delta_n^{2\alpha} \\ & \leq C \left(n^{-\frac{2\alpha}{2\alpha+d} + \xi} + \delta_n^{2\alpha} \right) \end{aligned}$$

with $\xi := \left(\frac{2\alpha(2\alpha+d+2)}{(2\alpha+d)((2\alpha+d)(1+p)+2p)} + \max\{\frac{d-2}{\beta d}, 0\} \right)p$ and constants $\tilde{c} > 0$ and $C > 0$ independent of n , Λ_n , ε_n , Γ_n , and δ_n . \blacksquare

In the end, we show Theorem 10 using Theorem 7 as well as Lemma 15.

Proof [Proof of Theorem 10] Let l be defined by $l := \lfloor \frac{n}{2} \rfloor + 1$, i.e. $l \geq \frac{n}{2}$. With this, Theorem 7 yields with probability \mathbb{P}^l not less than $1 - |\Lambda_n \times \Gamma_n|^{m_n} e^{-\tau}$ that

$$\mathcal{R}_{L_{J_T}, \mathbb{P}}(\hat{f}_{D_1, \lambda, \gamma}) - \mathcal{R}_{L_{J_T}, \mathbb{P}}^*$$

$$\begin{aligned}
&\leq \frac{c_1}{2} \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} + r_n^{2p} \left(\sum_{j=1}^{m_n} \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p l^{-1} + \tau l^{-1} \right) \\
&\leq c_1 \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} + r_n^{2p} \left(\sum_{j=1}^{m_n} \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p n^{-1} + \tau n^{-1} \right)
\end{aligned} \tag{37}$$

for all $(\lambda_j, \gamma_j) \in \Lambda_n \times \Gamma_n$, $j \in \{1, \dots, m_n\}$, simultaneously, where $c_1 > 0$ is a constant independent of n , τ , λ , and γ . Furthermore, the oracle inequality of (Steinwart and Christmann, 2008a, Theorem 7.2) for empirical risk minimization, $n - l \geq \frac{n}{2} - 1 \geq \frac{n}{4}$, and $\tau_n := \tau + \ln(1 + |\Lambda_n \times \Gamma_n|^{m_n})$ yield

$$\mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda_{D_2}, \gamma_{D_2}}) - \mathcal{R}_{L_{J_T}, P}^* \tag{38}$$

$$\begin{aligned}
&< 6 \left(\inf_{(\lambda_j, \gamma_j)_{j=1}^{m_n} \in (\Lambda_n \times \Gamma_n)^{m_n}} \mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda, \gamma}) - \mathcal{R}_{L_{J_T}, P}^* \right) + 512M^2 \frac{\tau_n}{n-l} \\
&< 6 \left(\inf_{(\lambda_j, \gamma_j)_{j=1}^{m_n} \in (\Lambda_n \times \Gamma_n)^{m_n}} \mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda, \gamma}) - \mathcal{R}_{L_{J_T}, P}^* \right) + 2048M^2 \frac{\tau_n}{n}
\end{aligned} \tag{39}$$

with probability P^{n-l} not less than $1 - e^{-\tau}$. With (37), (39) and Lemma 15 we can conclude

$$\begin{aligned}
&\mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda_{D_2}, \gamma_{D_2}}) - \mathcal{R}_{L_{J_T}, P}^* \\
&< 6 \left(\inf_{(\lambda_j, \gamma_j)_{j=1}^{m_n} \in (\Lambda_n \times \Gamma_n)^{m_n}} \mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda, \gamma}) - \mathcal{R}_{L_{J_T}, P}^* \right) + 2048M^2 \frac{\tau_n}{n} \\
&\leq 6c_1 \left(\inf_{(\lambda_j, \gamma_j)_{j=1}^{m_n} \in (\Lambda_n \times \Gamma_n)^{m_n}} \left(\sum_{j \in J_T} \lambda_j \gamma_j^{-d} + \left(\frac{\max_{j \in J_T} \gamma_j}{\min_{j \in J_T} \gamma_j} \right)^d \max_{j \in J_T} \gamma_j^{2\alpha} \right. \right. \\
&\quad \left. \left. + r_n^{2p} \left(\sum_{j=1}^{m_n} \lambda_j^{-1} \gamma_j^{-\frac{d+2p}{p}} P_X(A_j) \right)^p n^{-1} \right) + \tau n^{-1} \right) + 2048M^2 \frac{\tau_n}{n} \\
&\leq 6c_1 \left(C \left(n^{-\frac{2\alpha}{2\alpha+d} + \xi} + \delta_n^{2\alpha} \right) + \tau n^{-1} \right) + 2048M^2 \frac{\tau_n}{n} \\
&\leq 12c_1 C n^{-\frac{2\alpha}{2\alpha+d} + \xi} + (6c_1 \tau + 2048M^2 \tau_n) n^{-1}
\end{aligned}$$

with probability P^n not less than $1 - (1 + |\Lambda_n \times \Gamma_n|^{m_n}) e^{-\tau}$. Finally, a variable transformation yields

$$\begin{aligned}
&\mathcal{R}_{L_{J_T}, P}(\hat{f}_{D_1, \lambda_{D_2}, \gamma_{D_2}}) - \mathcal{R}_{L_{J_T}, P}^* \\
&< 12c_1 C n^{-\frac{2\alpha}{2\alpha+d} + \xi} + (6c_1 (\tau + \ln(1 + |\Lambda_n \times \Gamma_n|^{m_n})) \\
&\quad + 2048M^2 (\tau + 2 \ln(1 + |\Lambda_n \times \Gamma_n|^{m_n}))) n^{-1} \\
&\leq 12c_1 C n^{-\frac{2\alpha}{2\alpha+d} + \xi} + (6c_1 + 2048M^2) (\tau + 2m_n \ln(1 + |\Lambda_n \times \Gamma_n|)) n^{-1} \\
&\leq 12c_1 C n^{-\frac{2\alpha}{2\alpha+d} + \xi} + (6c_1 + 2048M^2) \left(\tau + 2 \cdot 8^d r_n^{-d} \ln(1 + |\Lambda_n \times \Gamma_n|) \right) n^{-1}
\end{aligned}$$

$$\begin{aligned}
&= 12c_1 C n^{-\frac{2\alpha}{2\alpha+d}+\xi} + (6c_1 + 2048M^2) \left(\tau n^{-1} + 2 \cdot 8^d c^{-d} \ln(1 + |\Lambda_n \times \Gamma_n|) n^{-\frac{\beta-1}{\beta}} \right) \\
&< \left(12c_1 C + (6c_1 + 2048M^2) \left(\tau + 2 \cdot 8^d c^{-d} \ln(1 + |\Lambda_n \times \Gamma_n|) \right) \right) n^{-\frac{2\alpha}{2\alpha+d}+\xi}
\end{aligned}$$

with probability P^n not less than $1 - e^{-\tau}$, where we used

$$\alpha < \frac{\beta-1}{2}d \iff n^{-\frac{\beta-1}{\beta}} < n^{-\frac{2\alpha}{2\alpha+d}}$$

in the last step. ■

References

- R. A. Adams and J. J. F. Fournier. *Sobolev Spaces*. Academic Press, New York, 2nd edition, 2003.
- N. Aronszajn. Theory of reproducing kernels. *Trans. Amer. Math. Soc.*, 68:337–404, 1950.
- K.P. Bennett and J. A. Blue. A support vector machine approach to decision trees. In *Neural Networks Proceedings, 1998. IEEE World Congress on Computational Intelligence. The 1998 IEEE International Joint Conference on*, volume 3, pages 2396–2401 vol.3, 1998. doi: 10.1109/IJCNN.1998.687237.
- H. Berens and R. DeVore. Quantitative Korovkin Theorems for Positive Linear Operators on L_p -Spaces. *Trans. Amer. Math. Soc.*, 245:pp. 349–361, 1978. ISSN 00029947. URL <http://www.jstor.org/stable/1998871>.
- A. Berlinet and C. Thomas-Agnan. *Reproducing Kernel Hilbert Spaces in Probability and Statistics*. Kluwer, Boston, 2004.
- E. Blanzieri and F. Melgani. Nearest Neighbor Classification of Remote Sensing Images With the Maximal Margin Principle. *Geoscience and Remote Sensing, IEEE Transactions on*, 46(6):1804–1811, 2008. ISSN 0196-2892. doi: 10.1109/TGRS.2008.916090.
- Enrico Blanzieri and Anton Bryl. Instance-based spam filtering using SVM nearest neighbor classifier. In *In Proceedings of FLAIRS 2007*, pages 441–442, 2007a.
- Enrico Blanzieri and Anton Bryl. Evaluation of the Highest Probability SVM Nearest Neighbor Classifier with Variable Relative Error Cost. In *In Proceedings of Fourth Conference on Email and Anti-Spam, CEAS’2007*, 2007b.
- Léon Bottou and Vladimir Vapnik. Local Learning Algorithms. *Neural Computation*, 4: 888–900, 1992.
- A. Caponnetto and E. De Vito. Optimal rates for regularized least squares algorithm. *Found. Comput. Math.*, 7:331–368, 2007.

- B. Carl and I. Stephani. *Entropy, Compactness, and the Approximation of Operators*. Cambridge Tracts in Mathematics. Cambridge University Press, Cambridge, 1990. ISBN 9780521330114. URL <http://books.google.de/books?id=x1BNHvgP7PIC>.
- Chih-Chung Chang and Chih-Jen Lin. LIBSVM: A library for support vector machines. *ACM Transactions on Intelligent Systems and Technology*, 2:27:1–27:27, 2011. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Fu Chang, Chien-Yang Guo, Xiao-Rong Lin, and Chi-Jen Lu. Tree Decomposition for Large-Scale SVM Problems. *J. Mach. Learn. Res.*, 11:2935–2972, December 2010. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1756006.1953027>.
- Haibin Cheng, Pang-Ning Tan, and Rong Jin. Localized Support Vector Machine and Its Efficient Algorithm. In *SIAM International Conference on Data Mining*, 2007.
- Haibin Cheng, Pang-Ning Tan, and Rong Jin. Efficient Algorithm for Localized Support Vector Machine. *Knowledge and Data Engineering, IEEE Transactions on*, 22(4):537–549, 2010. ISSN 1041-4347. doi: 10.1109/TKDE.2009.116.
- Ronan Collobert, Samy Bengio, and Yoshua Bengio. A Parallel Mixture of SVMs for Very Large Scale Problems. In *In Advances in Neural Information Processing Systems*, pages 633–640, 2001.
- F. Cucker and S. Smale. On the mathematical foundations of learning. *Bull. Amer. Math. Soc.*, 39:1–49, 2002.
- Sanjoy Dasgupta. Lecture 1: Clustering in metric spaces. CSE 291: Topics in unsupervised learning, 2008. URL <http://cseweb.ucsd.edu/~dasgupta/291/lec1.pdf>.
- E. De Vito, A. Caponnetto, and L. Rosasco. Model selection for regularized least-squares algorithm in learning theory. *Found. Comput. Math.*, 5:59–85, 2005.
- R.A. DeVore and V.A. Popov. Interpolation of Besov Spaces. *Trans. Amer. Math. Soc.*, 305(1):pp. 397–414, 1988.
- Ronald Alvin DeVore and George G. Lorentz. *Constructive Approximation*. Grundlehren der mathematischen Wissenschaften. Springer-Verlag, Berlin, 1993. ISBN 3-540-50627-6. URL http://books.google.de/books?id=cDqNW6k7_ZwC.
- Mona Eberts and Ingo Steinwart. Optimal learning rates for least squares SVMs using Gaussian kernels. In J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F.C.N. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 1539–1547. 2011.
- Mona Eberts and Ingo Steinwart. Optimal regression rates for SVMs using Gaussian kernels. *Electron. J. Statist.*, 7:1–42, 2013. ISSN 1935-7524. doi: 10.1214/12-EJS760.
- D. E. Edmunds and H. Triebel. *Function Spaces, Entropy Numbers, Differential Operators*. Cambridge University Press, Cambridge, 1996.

- Teofilo F. Gonzalez. Clustering to minimize the maximum intercluster distance. *Theoretical Computer Science*, 38(0):293–306, 1985. ISSN 0304-3975. doi: [http://dx.doi.org/10.1016/0304-3975\(85\)90224-5](http://dx.doi.org/10.1016/0304-3975(85)90224-5). URL <http://www.sciencedirect.com/science/article/pii/0304397585902245>.
- Hans Peter Graf, Eric Cosatto, Leon Bottou, Igor Durdanovic, and Vladimir Vapnik. Parallel Support Vector Machines: The Cascade SVM. In *In Advances in Neural Information Processing Systems*, pages 521–528. MIT Press, 2005.
- L. Györfi, M. Kohler, A. Krzyżak, and H. Walk. *A Distribution-Free Theory of Nonparametric Regression*. Springer-Verlag, New York, 2002.
- Robert Hable. Universal consistency of localized versions of regularized kernel methods. *J. Mach. Learn. Res.*, 14(1):153–186, January 2013. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=2502581.2502586>.
- Marius Kloft and Gilles Blanchard. On the convergence rate of ℓ_p -norm multiple kernel learning. *J. Mach. Learn. Res.*, 13:2465–2502, Aug 2012.
- R. Koenker. *Quantile Regression*. Econometric Society Monographs. Cambridge University Press, 1st edition, 2005. ISBN 9780521608275. URL <http://books.google.de/books?id=hdk7V4NXsgC>.
- S. Mendelson and J. Neeman. Regularization in kernel learning. *Ann. Statist.*, 38:526–565, 2010.
- Nicola Segata and Enrico Blanzieri. Empirical assessment of classification accuracy of local SVM. Technical report, University of Trento, Information Engineering and Computer Science, 2008. URL eprints.biblio.unitn.it/1398/1/014.pdf.
- Nicola Segata and Enrico Blanzieri. Fast and Scalable Local Kernel Machines. *J. Mach. Learn. Res.*, 11:1883–1926, August 2010. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1756006.1859915>.
- J. S. Simonoff. *Smoothing Methods in Statistics*. Springer, New York, 1996.
- S. Smale and D.-X. Zhou. Learning theory estimates via integral operators and their approximations. *Constr. Approx.*, 26:153–172, 2007.
- I. Steinwart and A. Christmann. *Support Vector Machines*. Springer, New York, 2008a.
- I. Steinwart and A. Christmann. How SVMs can estimate quantiles and the median. In J.C. Platt, D. Koller, Y. Singer, and S. Roweis, editors, *Advances in Neural Information Processing Systems 20*, pages 305–312, Cambridge, MA, 2008b. MIT Press.
- I. Steinwart and A. Christmann. Estimating conditional quantiles with the help of the pinball loss. *Bernoulli*, 17:211–225, 2011.
- I. Steinwart, D. Hush, and C. Scovel. Optimal rates for regularized least squares regression. In S. Dasgupta and A. Klivans, editors, *Proceedings of the 22nd Annual Conference on Learning Theory*, pages 79–93, 2009.

- I. Steinwart, D. Hush, and C. Scovel. Training SVMs without offset. *J. Mach. Learn. Res.*, 12:141–202, 2011.
- Taiji Suzuki. Unifying Framework for Fast Learning Rate of Non-Sparse Multiple Kernel Learning. In J. Shawe-Taylor, R.S. Zemel, P. Bartlett, F.C.N. Pereira, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems 24*, pages 1575–1583. 2011.
- Hans Triebel. *Theory of Function Spaces*. Birkhäuser, Basel, repr. of the 1983 edition, 2010. ISBN 978-3-0346-0415-4.
- Ivor W. Tsang, James T. Kwok, and Pak-Ming Cheung. Core Vector Machines: Fast SVM Training on Very Large Data Sets. *J. Mach. Learn. Res.*, 6:363–392, December 2005. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1046920.1058114>.
- Ivor W. Tsang, Andras Kocsor, and James T. Kwok. Simpler core vector machines with enclosing balls. In *Proceedings of the 24th international conference on Machine learning*, ICML '07, pages 911–918, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-793-3. doi: 10.1145/1273496.1273611. URL <http://doi.acm.org/10.1145/1273496.1273611>.
- Vladimir Vapnik and Léon Bottou. Local Algorithms for Pattern Recognition and Dependencies Estimation. *Neural Computation*, 5:893–909, 1993. doi: 10.1162/neco.1993.5.6.893.
- Donghui Wu, Kristin P. Bennett, Nello Cristianini, and John Shawe-Taylor. Large Margin Trees for Induction and Transduction. In *Proceedings of the Sixteenth International Conference on Machine Learning*, ICML '99, pages 474–483, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc. ISBN 1-55860-612-2. URL <http://dl.acm.org/citation.cfm?id=645528.657622>.
- Alon Zakai and Ya'acov Ritov. Consistency and Localizability. *J. Mach. Learn. Res.*, 10:827–856, June 2009. ISSN 1532-4435. URL <http://dl.acm.org/citation.cfm?id=1577069.1577099>.
- Hao Zhang, A.C. Berg, M. Maire, and J. Malik. SVM-KNN: Discriminative Nearest Neighbor Classification for Visual Category Recognition. In *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, volume 2, pages 2126–2136, 2006. doi: 10.1109/CVPR.2006.301.

Appendix A

For the sake of completeness, we present in the following some tables containing the computational results achieved by the LS-, VP-, and RC-SVMs for all real and artificial data set types. Here, the training and test times, given in seconds, are averaged over all successful runs. Moreover, for the test and L_2 -errors, we also stated the mean of all runs plus/minus the standard deviation. The same is true for the number of working sets (# of ws), except for the LS-SVMs, where we always have one working set by construction. The last two columns contain median, minimum, and maximum of the working set sizes appearing during the various runs.

	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	1.38	0.48	0.7142 \pm 0.0097	1.00		
	2 500 (10 000)	100	6.42	0.99	0.6323 \pm 0.0074	1.00		
	5 000 (10 000)	100	30.19	1.55	0.5707 \pm 0.0070	1.00		
	10 000 (50 000)	10	138.01	12.14	0.4909 \pm 0.0089	1.00		
	25 000 (50 000)	10	922.90	34.19	0.3816 \pm 0.0042	1.00		
	50 000 (50 000)	3	3788	176.68	0.3117 \pm 0.0012	1.00		
	100 000 (50 000)	3	16 353	507.38	0.2417 \pm 0.0102	1.00		
VP-SVM radius = 2	1 000 (10 000)	100	5.85	0.73	0.7521 \pm 0.0108	83.73 \pm 1.86	6	[1, 156]
	2 500 (10 000)	100	7.44	1.15	0.6602 \pm 0.0082	101.92 \pm 2.04	10	[1, 481]
	5 000 (10 000)	100	9.37	1.83	0.6011 \pm 0.0079	106.72 \pm 1.95	24	[1, 987]
	10 000 (50 000)	10	13.45	13.68	0.5238 \pm 0.0076	120.20 \pm 1.99	40	[1, 952]
	25 000 (50 000)	10	26.73	31.07	0.4134 \pm 0.0057	131.00 \pm 2.26	80	[1, 3204]
	50 000 (50 000)	3	57.41	183.71	0.3449 \pm 0.0047	139.67 \pm 4.16	155	[3, 4334]
	100 000 (50 000)	3	171.12	493.44	0.2658 \pm 0.0028	154.33 \pm 5.51	271	[3, 8121]
	250 000 (50 000)	3	1128	1169	0.1924 \pm 0.0016	166.33 \pm 3.79	533	[3, 22633]
	500 000 (50 000)	3	5349	3020	0.1608 \pm 0.0024	178.33 \pm 0.58	987	[2, 44585]
VP-SVM radius = 3	1 000 (10 000)	100	2.79	0.59	0.7262 \pm 0.0092	37.39 \pm 1.85	15	[1, 193]
	2 500 (10 000)	100	3.97	1.09	0.6379 \pm 0.0072	44.40 \pm 1.84	25	[1, 493]
	5 000 (10 000)	100	5.58	1.74	0.5845 \pm 0.0062	47.59 \pm 1.66	50	[1, 1008]
	10 000 (50 000)	10	10.64	13.43	0.5075 \pm 0.0042	51.20 \pm 1.62	85	[1, 1918]
	25 000 (50 000)	10	31.83	32.38	0.4026 \pm 0.0044	58.60 \pm 1.65	142	[3, 4971]
	50 000 (50 000)	3	120.10	208.34	0.3256 \pm 0.0021	61.33 \pm 2.08	274	[4, 10005]
	100 000 (50 000)	3	444.39	523.39	0.2539 \pm 0.0034	64.67 \pm 1.53	484	[18, 20082]
	250 000 (50 000)	3	2825	1274	0.1859 \pm 0.0018	68.33 \pm 1.15	1169	[21, 49558]
	500 000 (50 000)	3	12 882	3786	0.1540 \pm 0.0010	72.00 \pm 1.00	2366	[54, 74787]
VP-SVM radius = 4	1 000 (10 000)	100	0.87	0.53	0.7095 \pm 0.0102	4.40 \pm 0.49	165	[42, 489]
	2 500 (10 000)	100	2.14	1.05	0.6317 \pm 0.0062	5.02 \pm 0.49	382	[83, 1248]
	5 000 (10 000)	100	7.41	1.74	0.5737 \pm 0.0073	4.94 \pm 0.55	723	[229, 2401]
	10 000 (50 000)	10	28.58	13.62	0.5007 \pm 0.0070	5.20 \pm 0.42	1518	[488, 4514]
	25 000 (50 000)	10	201.17	33.16	0.3937 \pm 0.0049	5.70 \pm 0.67	2533	[566, 11113]
	50 000 (50 000)	3	827.41	218.85	0.3225 \pm 0.0028	5.33 \pm 0.58	7530	[2870, 24509]
	100 000 (50 000)	3	2823	548.96	0.2418 \pm 0.0007	7.00 \pm 0.00	6356	[2621, 40018]
	250 000 (50 000)	3	20 010	1434	0.1689 \pm 0.0041	6.33 \pm 0.58	17 721	[6061, 110464]
	500 000 (50 000)	0 (3)	NA	NA	NA \pm NA	7.33 \pm 1.53	68 182	[12469, 233675]
VP-SVM radius = 5	1 000 (10 000)	100	1.33	0.48	0.7138 \pm 0.0100	1.06 \pm 0.24	1000	[208, 1000]
	2 500 (10 000)	100	5.97	1.01	0.6326 \pm 0.0080	1.16 \pm 0.37	2500	[375, 2500]
	5 000 (10 000)	100	29.20	1.58	0.5705 \pm 0.0071	1.06 \pm 0.24	5000	[1810, 5000]
	10 000 (50 000)	10	131.83	12.48	0.4895 \pm 0.0066	1.10 \pm 0.32	10 000	[4338, 10000]
	25 000 (50 000)	10	832.08	33.09	0.3830 \pm 0.0045	1.20 \pm 0.42	25 000	[8524, 25000]
	50 000 (50 000)	3	3182	185.13	0.3151 \pm 0.0066	1.33 \pm 0.58	40 062	[19875, 50000]
	100 000 (50 000)	3	10 472	527.18	0.2427 \pm 0.0030	1.67 \pm 0.58	55 873	[42218, 100000]
	250 000 (50 000)	1 (3)	34 449	1445	0.1650 \pm 0.0000	3.00 \pm 0.00	86 655	[46661, 116684]
	500 000 (50 000)	0 (3)	NA	NA	NA \pm NA	1.33 \pm 0.58	375 000	[186510, 500000]

Table 3: LS- and VP-SVM results relating to the COVTYPE data sets

	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
RC-SVM # of ws = 1	1 000 (10 000)	100	1.76	0.30	0.7159 ± 0.0101	1.00 ± 0.00	1000	{1000}
	2 500 (10 000)	100	22.55	2.19	0.6324 ± 0.0073	1.00 ± 0.00	2500	{2500}
	5 000 (10 000)	100	97.30	9.16	0.5707 ± 0.0070	1.00 ± 0.00	5000	{5000}
	10 000 (50 000)	10	741.46	111.75	0.4909 ± 0.0089	1.00 ± 0.00	10 000	{10000}
	25 000 (50 000)	10	4629	291.75	0.3816 ± 0.0042	1.00 ± 0.00	25 000	{25000}
	50 000 (50 000)	3	11 416	521.84	0.3117 ± 0.0012	1.00 ± 0.00	50 000	{50000}
	100 000 (50 000)	3	19 921	622.99	0.2417 ± 0.0102	1.00 ± 0.00	100 000	{100000}
	250 000 (50 000)	0	NA	NA	NA \pm NA	NA \pm NA	NA	{NA}
	500 000 (50 000)	0	NA	NA	NA \pm NA	NA \pm NA	NA	{NA}
RC-SVM # of ws = 4	1 000 (10 000)	100	0.59	0.29	0.7527 ± 0.0143	4.00 ± 0.00	250	{250}
	2 500 (10 000)	100	4.76	2.56	0.6676 ± 0.0077	4.00 ± 0.00	625	{625}
	5 000 (10 000)	100	25.57	11.60	0.6254 ± 0.0074	4.00 ± 0.00	1250	{1250}
	10 000 (50 000)	10	89.55	83.44	0.5649 ± 0.0047	4.00 ± 0.00	2500	{2500}
	25 000 (50 000)	10	905.89	243.32	0.4746 ± 0.0062	4.00 ± 0.00	6250	{6250}
	50 000 (50 000)	3	2160	529.96	0.3990 ± 0.0032	4.00 ± 0.00	12 500	{12500}
	100 000 (50 000)	3	4191	826.40	0.3224 ± 0.0016	4.00 ± 0.00	25 000	{25000}
	250 000 (50 000)	3	29 909	1677	0.2251 ± 0.0012	4.00 ± 0.00	62 500	{62500}
	500 000 (50 000)	0 (3)	NA	NA	NA \pm NA	4.00 ± 0.00	125 000	{125000}
RC-SVM # of ws = 5	1 000 (10 000)	100	0.69	0.30	0.7674 ± 0.0142	5.00 ± 0.00	200	{200}
	2 500 (10 000)	100	3.59	2.97	0.6798 ± 0.0092	5.00 ± 0.00	500	{500}
	5 000 (10 000)	100	18.79	14.25	0.6404 ± 0.0082	5.00 ± 0.00	1000	{1000}
	10 000 (50 000)	10	72.89	112.97	0.5790 ± 0.0070	5.00 ± 0.00	2000	{2000}
	25 000 (50 000)	10	659.31	297.07	0.4886 ± 0.0048	5.00 ± 0.00	5000	{5000}
	50 000 (50 000)	3	1620	470.55	0.4216 ± 0.0024	5.00 ± 0.00	10 000	{10000}
	100 000 (50 000)	3	5441	1298	0.3420 ± 0.0040	5.00 ± 0.00	20 000	{20000}
	250 000 (50 000)	3	20 828	1614	0.2395 ± 0.0015	5.00 ± 0.00	50 000	{50000}
	500 000 (50 000)	3	81 660	2794	0.1931 ± 0.0007	5.00 ± 0.00	100 000	{100000}
RC-SVM # of ws = 6	1 000 (10 000)	100	0.73	0.29	0.7776 ± 0.0138	6.00 ± 0.00	167	[166 , 167]
	2 500 (10 000)	100	3.49	3.27	0.6911 ± 0.0101	6.00 ± 0.00	417	[416 , 417]
	5 000 (10 000)	100	15.61	16.57	0.6519 ± 0.0089	6.00 ± 0.00	833	[833 , 834]
	10 000 (50 000)	10	50.90	72.44	0.5937 ± 0.0060	6.00 ± 0.00	1667	[1666 , 1667]
	25 000 (50 000)	10	230.76	259.58	0.5035 ± 0.0055	6.00 ± 0.00	4167	[4166 , 4167]
	50 000 (50 000)	3	1406	425.29	0.4355 ± 0.0028	6.00 ± 0.00	8333	[8333 , 8334]
	100 000 (50 000)	3	5139	1203	0.3579 ± 0.0047	6.00 ± 0.00	16 667	[16666 , 16667]
	250 000 (50 000)	3	17 099	1577	0.2542 ± 0.0027	6.00 ± 0.00	41 667	[41666 , 41667]
	500 000 (50 000)	3	66 335	2534	0.2013 ± 0.0004	6.00 ± 0.00	83 333	[83333 , 83334]
RC-SVM # of ws = 7	1 000 (10 000)	100	0.72	0.30	0.7843 ± 0.0146	7.00 ± 0.00	143	[142 , 143]
	2 500 (10 000)	100	3.76	6.22	0.6991 ± 0.0090	7.00 ± 0.00	357	[357 , 358]
	5 000 (10 000)	100	9.51	10.49	0.6608 ± 0.0072	7.00 ± 0.00	714	[714 , 715]
	10 000 (50 000)	10	51.30	75.81	0.6045 ± 0.0078	7.00 ± 0.00	1429	[1428 , 1429]
	25 000 (50 000)	10	258.70	255.08	0.5163 ± 0.0057	7.00 ± 0.00	3571	[3571 , 3572]
	50 000 (50 000)	3	1087	440.26	0.4463 ± 0.0026	7.00 ± 0.00	7143	[7142 , 7143]
	100 000 (50 000)	3	3174	861.05	0.3765 ± 0.0033	7.00 ± 0.00	14 286	[14285 , 14286]
	250 000 (50 000)	3	15 472	1480	0.2638 ± 0.0012	7.00 ± 0.00	35 714	[35714 , 35715]
	500 000 (50 000)	3	56 601	2628	0.2094 ± 0.0013	7.00 ± 0.00	71 429	[71428 , 71429]

	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
RC-SVM # of ws = 10	1 000 (10 000)	100	0.99	0.32	0.8134 ± 0.0133	10.00 ± 0.00	100	{100}
	2 500 (10 000)	100	2.73	5.29	0.7188 ± 0.0090	10.00 ± 0.00	250	{250}
	5 000 (10 000)	100	6.61	15.24	0.6825 ± 0.0074	10.00 ± 0.00	500	{500}
	10 000 (50 000)	10	31.94	69.23	0.6278 ± 0.0079	10.00 ± 0.00	1000	{1000}
	25 000 (50 000)	10	259.18	304.44	0.5444 ± 0.0041	10.00 ± 0.00	2500	{2500}
	50 000 (50 000)	3	867.90	645.72	0.4904 ± 0.0044	10.00 ± 0.00	5000	{5000}
	100 000 (50 000)	3	2065	1204	0.4120 ± 0.0076	10.00 ± 0.00	10 000	{10000}
	250 000 (50 000)	3	9534	1398	0.3039 ± 0.0017	10.00 ± 0.00	25 000	{25000}
	500 000 (50 000)	3	43 146	2689	0.2300 ± 0.0008	10.00 ± 0.00	50 000	{50000}
RC-SVM # of ws = 20	1 000 (10 000)	100	1.67	0.33	0.8646 ± 0.0082	20.00 ± 0.00	50	{50}
	2 500 (10 000)	100	3.16	4.43	0.7574 ± 0.0076	20.00 ± 0.00	125	{125}
	5 000 (10 000)	100	3.99	13.46	0.7172 ± 0.0069	20.00 ± 0.00	250	{250}
	10 000 (50 000)	10	22.55	170.42	0.6770 ± 0.0068	20.00 ± 0.00	500	{500}
	25 000 (50 000)	10	87.06	236.67	0.5958 ± 0.0054	20.00 ± 0.00	1250	{1250}
	50 000 (50 000)	3	150.59	269.68	0.5470 ± 0.0029	20.00 ± 0.00	2500	{2500}
	100 000 (50 000)	3	885.28	685.50	0.4760 ± 0.0028	20.00 ± 0.00	5000	{5000}
	250 000 (50 000)	3	4559	1634	0.3752 ± 0.0057	20.00 ± 0.00	12 500	{12500}
	500 000 (50 000)	3	19 033	2802	0.2971 ± 0.0008	20.00 ± 0.00	25 000	{25000}
RC-SVM # of ws = 50	1 000 (10 000)	100	3.75	0.38	0.9280 ± 0.0046	50.00 ± 0.00	20	{20}
	2 500 (10 000)	100	6.58	5.52	0.8220 ± 0.0073	50.00 ± 0.00	50	{50}
	5 000 (10 000)	100	6.04	15.28	0.7679 ± 0.0081	50.00 ± 0.00	100	{100}
	10 000 (50 000)	10	11.50	124.14	0.7245 ± 0.0058	50.00 ± 0.00	200	{200}
	25 000 (50 000)	10	19.95	277.44	0.6551 ± 0.0056	50.00 ± 0.00	500	{500}
	50 000 (50 000)	3	54.19	221.31	0.6101 ± 0.0032	50.00 ± 0.00	1000	{1000}
	100 000 (50 000)	3	319.38	1018	0.5600 ± 0.0023	50.00 ± 0.00	2000	{2000}
	250 000 (50 000)	3	2680	1659	0.4688 ± 0.0037	50.00 ± 0.00	5000	{5000}
	500 000 (50 000)	3	6984	2352	0.4011 ± 0.0022	50.00 ± 0.00	10 000	{10000}
RC-SVM # of ws = 100	1 000 (10 000)	100	7.14	0.42	0.9562 ± 0.0035	100.00 ± 0.00	10	{10}
	2 500 (10 000)	100	12.05	5.85	0.8752 ± 0.0049	100.00 ± 0.00	25	{25}
	5 000 (10 000)	100	10.70	16.60	0.8239 ± 0.0061	100.00 ± 0.00	50	{50}
	10 000 (50 000)	10	12.05	92.09	0.7593 ± 0.0069	100.00 ± 0.00	100	{100}
	25 000 (50 000)	10	21.60	276.46	0.6926 ± 0.0037	100.00 ± 0.00	250	{250}
	50 000 (50 000)	3	32.39	206.27	0.6493 ± 0.0033	100.00 ± 0.00	500	{500}
	100 000 (50 000)	3	141.39	798.77	0.6082 ± 0.0022	100.00 ± 0.00	1000	{1000}
	250 000 (50 000)	3	1325	1970	0.5306 ± 0.0029	100.00 ± 0.00	2500	{2500}
	500 000 (50 000)	3	3053	2315	0.4670 ± 0.0025	100.00 ± 0.00	5000	{5000}
RC-SVM # of ws = 150	1 000 (10 000)	100	11.98	0.63	0.9650 ± 0.0023	150.00 ± 0.00	7	[6 , 7]
	2 500 (10 000)	100	15.21	4.66	0.9086 ± 0.0040	150.00 ± 0.00	17	[16 , 17]
	5 000 (10 000)	100	16.94	14.18	0.8610 ± 0.0053	150.00 ± 0.00	33	[33 , 34]
	10 000 (50 000)	10	18.06	142.93	0.7936 ± 0.0068	150.00 ± 0.00	67	[66 , 67]
	25 000 (50 000)	10	16.47	183.81	0.7141 ± 0.0039	150.00 ± 0.00	167	[166 , 167]
	50 000 (50 000)	3	26.25	187.73	0.6742 ± 0.0008	150.00 ± 0.00	333	[333 , 334]
	100 000 (50 000)	3	100.66	755.96	0.6337 ± 0.0020	150.00 ± 0.00	667	[666 , 667]
	250 000 (50 000)	3	537.02	1455	0.5658 ± 0.0004	150.00 ± 0.00	1667	[1666 , 1667]
	500 000 (50 000)	3	1859	2336	0.5081 ± 0.0019	150.00 ± 0.00	3333	[3333 , 3334]

Table 4: RC-SVM results relating to the COVTYPE data sets

	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	1.21	0.25	0.1916 ±0.0055	1.00		
	2 500 (10 000)	100	8.25	0.57	0.1842 ±0.0043	1.00		
	5 000 (10 000)	100	33.98	1.03	0.1672 ±0.0032	1.00		
	10 000 (50 000)	10	129.26	11.76	0.1596 ±0.0032	1.00		
	25 000 (50 000)	10	791.14	52.75	0.1453 ±0.0020	1.00		
	50 000 (50 000)	3	3029	169.78	0.1410 ±0.0022	1.00		
	100 000 (50 000)	3	11 078	201.54	0.1302 ±0.0008	1.00		
	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
VP-SVM radius = 1	1 000 (10 000)	100	1.81	0.46	0.2586 ±0.0090	21.62 ±1.23	15	[1, 469]
	2 500 (10 000)	100	2.90	0.90	0.2243 ±0.0058	25.03 ±1.58	28	[1, 1289]
	5 000 (10 000)	100	5.19	1.34	0.1891 ±0.0036	26.50 ±1.63	62	[2, 2520]
	10 000 (50 000)	10	12.27	9.21	0.1724 ±0.0033	28.50 ±1.51	84	[1, 5018]
	25 000 (50 000)	10	43.56	21.51	0.1522 ±0.0024	34.10 ±1.52	209	[1, 6015]
	50 000 (50 000)	3	245.63	169.75	0.1462 ±0.0027	39.67 ±2.31	275	[3, 18273]
	100 000 (50 000)	3	932.30	347.92	0.1352 ±0.0024	41.00 ±1.00	535	[6, 30230]
VP-SVM radius = 2	250 000 (50 000)	3	9404	869.18	0.1277 ±0.0014	44.67 ±0.58	1120	[7, 112012]
	400 000 (50 000)	3	17 110	1437	0.1152 ±0.0014	44.67 ±1.15	1804	[10, 116673]
	1 000 (10 000)	100	1.02	0.40	0.2056 ±0.0070	4.32 ±0.72	86	[5, 956]
	2 500 (10 000)	100	2.69	0.81	0.1953 ±0.0050	4.84 ±0.72	200	[8, 2263]
	5 000 (10 000)	100	10.51	1.25	0.1743 ±0.0040	5.06 ±0.81	304	[3, 4670]
	10 000 (50 000)	10	45.12	8.41	0.1636 ±0.0022	4.70 ±0.48	583	[63, 8697]
	25 000 (50 000)	10	264.42	34.60	0.1466 ±0.0038	5.70 ±0.67	1360	[51, 19740]
VP-SVM radius = 3	50 000 (50 000)	3	1268	186.33	0.1445 ±0.0008	6.33 ±1.15	1663	[180, 42173]
	100 000 (50 000)	3	3670	425.33	0.1279 ±0.0016	6.33 ±0.58	7167	[359, 65686]
	250 000 (50 000)	0 (3)	NA	NA	NA ±NA	7.33 ±2.08	34 091	[547, 210793]
	400 000 (50 000)	0 (3)	NA	NA	NA ±NA	7.33 ±1.53	54 545	[738, 336202]
	1 000 (10 000)	100	1.18	0.27	0.1932 ±0.0061	1.41 ±0.57	957	[36, 1000]
	2 500 (10 000)	100	5.28	0.59	0.1875 ±0.0055	1.60 ±0.53	2371	[55, 2500]
	5 000 (10 000)	100	24.31	1.08	0.1687 ±0.0035	1.36 ±0.48	4795	[188, 5000]
VP-SVM radius = 4	10 000 (50 000)	10	116.00	8.76	0.1592 ±0.0032	1.10 ±0.32	10 000	[615, 10000]
	25 000 (50 000)	10	580.23	41.56	0.1473 ±0.0029	1.80 ±0.42	21 051	[752, 25000]
	50 000 (50 000)	3	3130	186.34	0.1429 ±0.0015	1.67 ±0.58	42 736	[1861, 50000]
	100 000 (50 000)	3	8886	371.73	0.1303 ±0.0029	2.33 ±0.58	12 300	[2023, 97977]
	250 000 (50 000)	0 (3)	NA	NA	NA ±NA	2.33 ±0.58	107 143	[2796, 247204]
	400 000 (50 000)	0 (3)	NA	NA	NA ±NA	2.00 ±0.00	200 000	[4773, 395227]
	VP-SVM radius = 5	1 000 (10 000)	100	1.15	0.23	0.1913 ±0.0056	1.00 ±0.00	1000
2 500 (10 000)		100	6.45	0.57	0.1842 ±0.0043	1.00 ±0.00	2500	{2500}
5 000 (10 000)		100	28.18	1.06	0.1672 ±0.0032	1.00 ±0.00	5000	{5000}
10 000 (50 000)		10	119.95	8.58	0.1596 ±0.0032	1.00 ±0.00	10 000	{10000}
25 000 (50 000)		10	807.92	72.42	0.1452 ±0.0020	1.10 ±0.32	25 000	{6292, 25000}
50 000 (50 000)		3	3539	197.21	0.1410 ±0.0022	1.00 ±0.00	50 000	{50000}
100 000 (50 000)		3	12 679	389.43	0.1302 ±0.0008	1.00 ±0.00	100 000	{100000}
VP-SVM radius = 6	250 000 (50 000)	0 (3)	NA	NA	NA ±NA	1.00 ±0.00	250 000	{250000}
	400 000 (50 000)	0 (3)	NA	NA	NA ±NA	1.00 ±0.00	400 000	{400000}
	1 000 (10 000)	100	1.02	0.45	0.2341 ±0.0048	10.00 ±0.00	100	{100}
	2 500 (10 000)	100	1.45	0.51	0.2136 ±0.0030	10.00 ±0.00	250	{250}
	5 000 (10 000)	100	3.42	1.30	0.1928 ±0.0026	10.00 ±0.00	500	{500}
	10 000 (50 000)	10	12.03	13.35	0.1831 ±0.0028	10.00 ±0.00	1000	{1000}
	25 000 (50 000)	10	58.63	22.87	0.1650 ±0.0013	10.00 ±0.00	2500	{2500}
RC-SVM # of ws = 1	50 000 (50 000)	3	264.45	47.21	0.1613 ±0.0008	10.00 ±0.00	5000	{5000}
	100 000 (50 000)	3	1166	190.86	0.1450 ±0.0022	10.00 ±0.00	10 000	{10000}
	250 000 (50 000)	3	7119	643.38	0.1374 ±0.0012	10.00 ±0.00	25 000	{25000}
	400 000 (50 000)	3	18 495	1317	0.1306 ±0.0014	10.00 ±0.00	40 000	{40000}
	1 000 (10 000)	100	1.61	0.50	0.2736 ±0.0071	20.00 ±0.00	50	{50}
	2 500 (10 000)	100	1.76	0.55	0.2365 ±0.0038	20.00 ±0.00	125	{125}
	5 000 (10 000)	100	2.84	1.33	0.2067 ±0.0024	20.00 ±0.00	250	{250}
RC-SVM # of ws = 2	10 000 (50 000)	10	6.79	14.07	0.1946 ±0.0026	20.00 ±0.00	500	{500}
	25 000 (50 000)	10	31.20	23.20	0.1730 ±0.0021	20.00 ±0.00	1250	{1250}
	50 000 (50 000)	3	115.38	46.53	0.1699 ±0.0003	20.00 ±0.00	2500	{2500}
	100 000 (50 000)	3	524.48	201.80	0.1548 ±0.0009	20.00 ±0.00	5000	{5000}
	250 000 (50 000)	3	3688	616.18	0.1418 ±0.0006	20.00 ±0.00	12 500	{12500}
	400 000 (50 000)	3	9783	1119	0.1374 ±0.0003	20.00 ±0.00	20 000	{20000}
	RC-SVM # of ws = 3	1 000 (10 000)	100	3.35	0.51	0.4617 ±0.0284	50.00 ±0.00	20
2 500 (10 000)		100	3.31	0.63	0.2958 ±0.0079	50.00 ±0.00	50	{50}
5 000 (10 000)		100	4.79	1.56	0.2351 ±0.0034	50.00 ±0.00	100	{100}
10 000 (50 000)		10	5.77	15.00	0.2125 ±0.0014	50.00 ±0.00	200	{200}
25 000 (50 000)		10	25.22	23.75	0.1867 ±0.0013	50.00 ±0.00	500	{500}
50 000 (50 000)		3	52.65	43.21	0.1812 ±0.0004	50.00 ±0.00	1000	{1000}
100 000 (50 000)		3	184.55	170.12	0.1666 ±0.0002	50.00 ±0.00	2000	{2000}
RC-SVM # of ws = 4	250 000 (50 000)	3	1324	588.68	0.1547 ±0.0003	50.00 ±0.00	5000	{5000}
	400 000 (50 000)	3	3825	1008	0.1491 ±0.0009	50.00 ±0.00	8000	{8000}

	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	1.13	0.41	0.1582 \pm 0.0052	1.00		
	2 500 (10 000)	100	4.08	0.80	0.1099 \pm 0.0028	1.00		
	5 000 (10 000)	100	18.37	1.31	0.1010 \pm 0.0016	1.00		
	10 000 (50 000)	10	82.65	9.29	0.0728 \pm 0.0006	1.00		
	25 000 (50 000)	10	526.48	24.21	0.0542 \pm 0.0004	1.00		
	50 000 (50 000)	3	2146	46.89	0.0448 \pm 0.0003	1.00		
	100 000 (50 000)	3	8907	246.99	0.0365 \pm 0.0002	1.00		
	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
VP-SVM radius = 1	1 000 (10 000)	100	4.22	0.58	0.2496 \pm 0.0074	58.07 \pm 1.99	10	[1, 83]
	2 500 (10 000)	100	5.79	0.95	0.1489 \pm 0.0043	76.80 \pm 1.93	17	[1, 177]
	5 000 (10 000)	100	7.72	1.40	0.1218 \pm 0.0029	93.44 \pm 2.24	21	[1, 336]
	10 000 (50 000)	10	11.40	9.93	0.0834 \pm 0.0023	115.70 \pm 2.63	29	[1, 618]
	25 000 (50 000)	10	24.17	30.30	0.0614 \pm 0.0014	154.10 \pm 3.35	43	[1, 1835]
	50 000 (50 000)	3	86.99	201.41	0.0455 \pm 0.0012	162.33 \pm 4.04	92	[2, 2808]
	100 000 (50 000)	3	264.22	356.52	0.0365 \pm 0.0019	199.00 \pm 7.81	127	[1, 5386]
VP-SVM radius = 2	1 000 (10 000)	100	1.48	0.47	0.2038 \pm 0.0061	18.08 \pm 0.69	67	[12, 110]
	2 500 (10 000)	100	1.92	0.89	0.1185 \pm 0.0034	18.93 \pm 0.74	169	[23, 276]
	5 000 (10 000)	100	3.42	1.42	0.0977 \pm 0.0024	19.41 \pm 0.57	339	[16, 539]
	10 000 (50 000)	10	7.24	14.77	0.0660 \pm 0.0015	20.40 \pm 0.52	348	[41, 1010]
	25 000 (50 000)	10	28.04	34.35	0.0513 \pm 0.0010	22.20 \pm 0.92	664	[64, 2516]
	50 000 (50 000)	3	196.07	235.05	0.0408 \pm 0.0002	23.33 \pm 0.58	1236	[171, 4336]
	100 000 (50 000)	3	688.80	410.05	0.0343 \pm 0.0006	24.33 \pm 0.58	2141	[306, 9959]
VP-SVM radius = 3	1 000 (10 000)	100	0.97	0.47	0.1857 \pm 0.0050	9.98 \pm 0.14	103	[84, 152]
	2 500 (10 000)	100	1.34	0.90	0.1109 \pm 0.0036	10.00 \pm 0.00	252	[226, 276]
	5 000 (10 000)	100	3.38	1.93	0.0966 \pm 0.0023	10.00 \pm 0.00	498	[467, 539]
	10 000 (50 000)	10	8.67	14.17	0.0660 \pm 0.0008	10.00 \pm 0.00	1000	[962, 1059]
	25 000 (50 000)	10	64.97	74.29	0.0507 \pm 0.0005	10.00 \pm 0.00	2494	[2447, 2566]
	50 000 (50 000)	3	310.74	248.32	0.0409 \pm 0.0004	10.00 \pm 0.00	5002	[4902, 5147]
	100 000 (50 000)	3	1127	421.73	0.0342 \pm 0.0004	10.00 \pm 0.00	9988	[9912, 10070]
VP-SVM radius = 4	1 000 (10 000)	100	1.12	0.39	0.1582 \pm 0.0052	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	3.94	0.81	0.1099 \pm 0.0028	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	21.67	1.36	0.1010 \pm 0.0016	1.00 \pm 0.00	5000	{5000}
	10 000 (50 000)	10	92.45	11.62	0.0728 \pm 0.0006	1.00 \pm 0.00	10 000	{10000}
	25 000 (50 000)	10	1052	110.79	0.0542 \pm 0.0004	1.00 \pm 0.00	25 000	{25000}
	50 000 (50 000)	3	3621	201.31	0.0448 \pm 0.0003	1.00 \pm 0.00	50 000	{50000}
	100 000 (50 000)	3	12 439	333.57	0.0365 \pm 0.0002	1.00 \pm 0.00	100 000	{100000}

	data set sizes	runs	train time	test time	test error	# of ws	ws size: median	ws size: range
RC-SVM # of ws = 1	1 000 (10 000)	100	1.08	0.38	0.1582 \pm 0.0052	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	4.01	0.85	0.1099 \pm 0.0028	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	17.47	1.29	0.1010 \pm 0.0016	1.00 \pm 0.00	5000	{5000}
	10 000 (50 000)	10	79.67	8.62	0.0728 \pm 0.0006	1.00 \pm 0.00	10 000	{10000}
	25 000 (50 000)	10	500.55	22.20	0.0542 \pm 0.0004	1.00 \pm 0.00	25 000	{25000}
	50 000 (50 000)	3	2388	116.48	0.0448 \pm 0.0003	1.00 \pm 0.00	50 000	{50000}
	100 000 (50 000)	3	8592	231.08	0.0365 \pm 0.0002	1.00 \pm 0.00	100 000	{100000}
RC-SVM # of ws = 5	1 000 (10 000)	100	0.66	0.39	0.2250 \pm 0.0075	5.00 \pm 0.00	200	{200}
	2 500 (10 000)	100	1.26	0.84	0.1698 \pm 0.0050	5.00 \pm 0.00	500	{500}
	5 000 (10 000)	100	3.73	1.45	0.1480 \pm 0.0031	5.00 \pm 0.00	1000	{1000}
	10 000 (50 000)	10	12.55	11.46	0.1093 \pm 0.0015	5.00 \pm 0.00	2000	{2000}
	25 000 (50 000)	10	88.31	27.18	0.0805 \pm 0.0009	5.00 \pm 0.00	5000	{5000}
	50 000 (50 000)	3	612.93	226.31	0.0624 \pm 0.0003	5.00 \pm 0.00	10 000	{10000}
	100 000 (50 000)	3	1623	301.96	0.0488 \pm 0.0001	5.00 \pm 0.00	20 000	{20000}
RC-SVM # of ws = 10	1 000 (10 000)	100	0.94	0.36	0.2648 \pm 0.0059	10.00 \pm 0.00	100	{100}
	2 500 (10 000)	100	1.40	0.90	0.2064 \pm 0.0049	10.00 \pm 0.00	250	{250}
	5 000 (10 000)	100	3.20	1.39	0.1806 \pm 0.0039	10.00 \pm 0.00	500	{500}
	10 000 (50 000)	10	6.99	10.92	0.1372 \pm 0.0027	10.00 \pm 0.00	1000	{1000}
	25 000 (50 000)	10	40.29	29.37	0.0992 \pm 0.0007	10.00 \pm 0.00	2500	{2500}
	50 000 (50 000)	3	199.33	168.55	0.0763 \pm 0.0006	10.00 \pm 0.00	5000	{5000}
	100 000 (50 000)	3	892.77	499.74	0.0590 \pm 0.0004	10.00 \pm 0.00	10 000	{10000}
RC-SVM # of ws = 20	1 000 (10 000)	100	1.61	0.37	0.3066 \pm 0.0038	20.00 \pm 0.00	50	{50}
	2 500 (10 000)	100	1.95	0.86	0.2461 \pm 0.0039	20.00 \pm 0.00	125	{125}
	5 000 (10 000)	100	2.95	1.35	0.2156 \pm 0.0043	20.00 \pm 0.00	250	{250}
	10 000 (50 000)	10	4.67	10.56	0.1704 \pm 0.0031	20.00 \pm 0.00	500	{500}
	25 000 (50 000)	10	22.59	29.35	0.1275 \pm 0.0008	20.00 \pm 0.00	1250	{1250}
	50 000 (50 000)	3	92.33	168.15	0.0964 \pm 0.0011	20.00 \pm 0.00	2500	{2500}
	100 000 (50 000)	3	1081	613.04	0.0731 \pm 0.0004	20.00 \pm 0.00	5000	{5000}
RC-SVM # of ws = 50	1 000 (10 000)	100	4.14	0.43	0.3478 \pm 0.0048	50.00 \pm 0.00	20	{20}
	2 500 (10 000)	100	3.74	0.84	0.2966 \pm 0.0026	50.00 \pm 0.00	50	{50}
	5 000 (10 000)	100	4.97	1.31	0.2694 \pm 0.0036	50.00 \pm 0.00	100	{100}
	10 000 (50 000)	10	5.46	10.90	0.2220 \pm 0.0029	50.00 \pm 0.00	200	{200}
	25 000 (50 000)	10	18.23	25.04	0.1715 \pm 0.0031	50.00 \pm 0.00	500	{500}
	50 000 (50 000)	3	39.18	143.68	0.1366 \pm 0.0014	50.00 \pm 0.00	1000	{1000}
	100 000 (50 000)	3	133.67	264.38	0.1036 \pm 0.0028	50.00 \pm 0.00	2000	{2000}
RC-SVM # of ws = 100	1 000 (10 000)	100	6.86	0.51	0.3502 \pm 0.0040	100.00 \pm 0.00	10	{10}
	2 500 (10 000)	100	7.37	0.92	0.3268 \pm 0.0031	100.00 \pm 0.00	25	{25}
	5 000 (10 000)	100	7.58	1.31	0.3123 \pm 0.0022	100.00 \pm 0.00	50	{50}
	10 000 (50 000)	10	8.58	8.93	0.2678 \pm 0.0032	100.00 \pm 0.00	100	{100}
	25 000 (50 000)	10	15.91	21.95	0.2077 \pm 0.0021	100.00 \pm 0.00	250	{250}
	50 000 (50 000)	3	27.94	132.07	0.1679 \pm 0.0023	100.00 \pm 0.00	500	{500}
	100 000 (50 000)	3	79.05	227.69	0.1341 \pm 0.0017	100.00 \pm 0.00	1000	{1000}
RC-SVM # of ws = 150	1 000 (10 000)	100	10.52	0.62	0.3468 \pm 0.0026	150.00 \pm 0.00	7	{6, 7}
	2 500 (10 000)	100	10.43	1.00	0.3328 \pm 0.0024	150.00 \pm 0.00	17	{16, 17}
	5 000 (10 000)	100	9.34	1.23	0.3395 \pm 0.0082	150.00 \pm 0.00	33	{33, 34}
	10 000 (50 000)	10	11.38	9.10	0.2943 \pm 0.0014	150.00 \pm 0.00	67	{66, 67}
	25 000 (50 000)	10	17.25	19.51	0.2336 \pm 0.0041	150.00 \pm 0.00	167	{166, 167}
	50 000 (50 000)	3	25.66	122.27	0.1927 \pm 0.0005	150.00 \pm 0.00	333	{333, 334}
	100 000 (50 000)	3	73.12	188.29	0.1538 \pm 0.0012	150.00 \pm 0.00	667	{666, 667}
RC-SVM # of ws = 200	1 000 (10 000)	100	13.70	0.64	0.3447 \pm 0.0021	200.00 \pm 0.00	5	{5}
	2 500 (10 000)	100	14.45	1.11	0.3346 \pm 0.0020	200.00 \pm 0.00	12	{12, 13}
	5 000 (10 000)	100	12.32	1.22	0.3534 \pm 0.0092	200.00 \pm 0.00	25	{25}
	10 000 (50 000)	10	13.19	8.93	0.3091 \pm 0.0024	200.00 \pm 0.00	50	{50}
	25 000 (50 000)	10	19.80	18.41	0.2512 \pm 0.0030	200.00 \pm 0.00	125	{125}
	50 000 (50 000)	3	26.89	87.32	0.2095 \pm 0.0033	200.00 \pm 0.00	250	{250}
	100 000 (50 000)	3	70.20	185.72	0.1664 \pm 0.0021	200.00 \pm 0.00	500	{500}

Table 6: Experimental results relating to the IJCNN1 data sets

	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	0.33	0.07	0.0284 \pm 0.0003	0.0541 \pm 0.0033	1.00		
	2 500 (10 000)	100	1.74	0.07	0.0275 \pm 0.0002	0.0461 \pm 0.0019	1.00		
	5 000 (10 000)	100	7.14	0.14	0.0269 \pm 0.0002	0.0396 \pm 0.0019	1.00		
	10 000 (10 000)	100	27.90	0.22	0.0265 \pm 0.0001	0.0323 \pm 0.0014	1.00		
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
VP-SVM $r = 0.1$	1 000 (10 000)	100	1.84	0.28	0.0287 \pm 0.0005	0.0561 \pm 0.0041	15.70 \pm 2.43	64	[27, 109]
	2 500 (10 000)	100	1.90	0.34	0.0267 \pm 0.0002	0.0368 \pm 0.0027	15.77 \pm 2.34	159	[59, 268]
	5 000 (10 000)	100	2.01	0.40	0.0260 \pm 0.0001	0.0248 \pm 0.0018	16.04 \pm 2.45	312	[125, 532]
	10 000 (10 000)	100	3.05	0.55	0.0257 \pm 0.0001	0.0197 \pm 0.0010	15.55 \pm 2.43	643	[245, 1054]
VP-SVM $r = 0.25$	1 000 (10 000)	100	0.65	0.16	0.0277 \pm 0.0004	0.0465 \pm 0.0037	6.17 \pm 0.88	162	[66, 251]
	2 500 (10 000)	100	0.75	0.18	0.0265 \pm 0.0002	0.0338 \pm 0.0024	6.12 \pm 0.87	408	[152, 645]
	5 000 (10 000)	100	1.58	0.24	0.0260 \pm 0.0001	0.0256 \pm 0.0016	6.17 \pm 0.79	810	[305, 1291]
	10 000 (10 000)	100	4.93	0.38	0.0258 \pm 0.0001	0.0213 \pm 0.0013	6.27 \pm 0.83	1595	[600, 2550]
VP-SVM $r = 0.5$	1 000 (10 000)	100	0.37	0.11	0.0275 \pm 0.0003	0.0447 \pm 0.0032	3.49 \pm 0.50	287	[121, 502]
	2 500 (10 000)	100	0.77	0.14	0.0266 \pm 0.0002	0.0355 \pm 0.0025	3.44 \pm 0.50	727	[329, 1202]
	5 000 (10 000)	100	2.31	0.19	0.0262 \pm 0.0001	0.0290 \pm 0.0021	3.41 \pm 0.49	1466	[627, 2511]
	10 000 (10 000)	100	8.88	0.31	0.0259 \pm 0.0001	0.0242 \pm 0.0019	3.51 \pm 0.50	2849	[1253, 5003]
VP-SVM $r = 1$	1 000 (10 000)	100	0.27	0.08	0.0282 \pm 0.0003	0.0517 \pm 0.0030	2.00 \pm 0.00	500	[266, 734]
	2 500 (10 000)	100	0.97	0.11	0.0272 \pm 0.0002	0.0426 \pm 0.0024	2.00 \pm 0.00	1250	[658, 1842]
	5 000 (10 000)	100	3.68	0.15	0.0268 \pm 0.0002	0.0382 \pm 0.0021	2.00 \pm 0.00	2500	[1224, 3776]
	10 000 (10 000)	100	14.93	0.44	0.0263 \pm 0.0001	0.0309 \pm 0.0020	2.00 \pm 0.00	5000	[2501, 7499]
VP-SVM $r = 2$	1 000 (10 000)	100	0.32	0.07	0.0284 \pm 0.0004	0.0541 \pm 0.0033	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	1.75	0.08	0.0275 \pm 0.0002	0.0461 \pm 0.0020	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	7.16	0.15	0.0269 \pm 0.0002	0.0396 \pm 0.0019	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	27.91	0.24	0.0265 \pm 0.0001	0.0323 \pm 0.0014	1.00 \pm 0.00	10 000	{10000}
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
RC-SVM # ws = 1	1 000 (10 000)	100	0.32	0.07	0.0284 \pm 0.0003	0.0540 \pm 0.0031	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	1.76	0.08	0.0275 \pm 0.0002	0.0461 \pm 0.0020	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	7.16	0.15	0.0269 \pm 0.0002	0.0396 \pm 0.0019	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	28.05	0.24	0.0265 \pm 0.0001	0.0323 \pm 0.0014	1.00 \pm 0.00	10 000	{10000}
RC-SVM # ws = 2	1 000 (10 000)	100	0.26	0.09	0.0287 \pm 0.0005	0.0571 \pm 0.0042	2.00 \pm 0.00	500	{500}
	2 500 (10 000)	100	0.94	0.11	0.0276 \pm 0.0002	0.0475 \pm 0.0018	2.00 \pm 0.00	1250	{1250}
	5 000 (10 000)	100	3.46	0.14	0.0273 \pm 0.0001	0.0444 \pm 0.0014	2.00 \pm 0.00	2500	{2500}
	10 000 (10 000)	100	14.37	0.28	0.0268 \pm 0.0001	0.0376 \pm 0.0016	2.00 \pm 0.00	5000	{5000}
RC-SVM # ws = 3	1 000 (10 000)	100	0.31	0.11	0.0288 \pm 0.0004	0.0583 \pm 0.0037	3.00 \pm 0.00	333	[333, 334]
	2 500 (10 000)	100	0.76	0.14	0.0277 \pm 0.0002	0.0488 \pm 0.0022	3.00 \pm 0.00	833	[833, 834]
	5 000 (10 000)	100	2.32	0.17	0.0274 \pm 0.0001	0.0456 \pm 0.0014	3.00 \pm 0.00	1667	[1666, 1667]
	10 000 (10 000)	100	9.58	0.27	0.0271 \pm 0.0001	0.0417 \pm 0.0015	3.00 \pm 0.00	3333	[3333, 3334]
RC-SVM # ws = 4	1 000 (10 000)	100	0.43	0.12	0.0289 \pm 0.0004	0.0588 \pm 0.0037	4.00 \pm 0.00	250	{250}
	2 500 (10 000)	100	0.69	0.16	0.0278 \pm 0.0003	0.0500 \pm 0.0025	4.00 \pm 0.00	625	{625}
	5 000 (10 000)	100	1.90	0.20	0.0275 \pm 0.0002	0.0469 \pm 0.0016	4.00 \pm 0.00	1250	{1250}
	10 000 (10 000)	100	6.95	0.27	0.0272 \pm 0.0001	0.0432 \pm 0.0010	4.00 \pm 0.00	2500	{2500}
RC-SVM # ws = 5	1 000 (10 000)	100	0.52	0.15	0.0287 \pm 0.0004	0.0577 \pm 0.0033	5.00 \pm 0.00	200	{200}
	2 500 (10 000)	100	0.67	0.18	0.0280 \pm 0.0002	0.0515 \pm 0.0022	5.00 \pm 0.00	500	{500}
	5 000 (10 000)	100	1.65	0.23	0.0276 \pm 0.0002	0.0480 \pm 0.0020	5.00 \pm 0.00	1000	{1000}
	10 000 (10 000)	100	5.52	0.30	0.0273 \pm 0.0001	0.0438 \pm 0.0012	5.00 \pm 0.00	2000	{2000}
RC-SVM # ws = 6	1 000 (10 000)	100	0.64	0.16	0.0288 \pm 0.0004	0.0586 \pm 0.0031	6.00 \pm 0.00	167	[166, 167]
	2 500 (10 000)	100	0.69	0.19	0.0281 \pm 0.0003	0.0525 \pm 0.0028	6.00 \pm 0.00	417	[416, 417]
	5 000 (10 000)	100	1.52	0.25	0.0277 \pm 0.0002	0.0490 \pm 0.0022	6.00 \pm 0.00	833	[833, 834]
	10 000 (10 000)	100	4.67	0.32	0.0273 \pm 0.0001	0.0446 \pm 0.0011	6.00 \pm 0.00	1667	[1666, 1667]
RC-SVM # ws = 10	1 000 (10 000)	100	1.05	0.22	0.0288 \pm 0.0003	0.0584 \pm 0.0029	10.00 \pm 0.00	100	{100}
	2 500 (10 000)	100	1.08	0.27	0.0282 \pm 0.0003	0.0534 \pm 0.0026	10.00 \pm 0.00	250	{250}
	5 000 (10 000)	100	1.39	0.34	0.0280 \pm 0.0002	0.0522 \pm 0.0023	10.00 \pm 0.00	500	{500}
	10 000 (10 000)	100	3.45	0.44	0.0275 \pm 0.0001	0.0471 \pm 0.0012	10.00 \pm 0.00	1000	{1000}
RC-SVM # ws = 20	1 000 (10 000)	100	2.28	0.34	0.0293 \pm 0.0004	0.0626 \pm 0.0027	20.00 \pm 0.00	50	{50}
	2 500 (10 000)	100	2.06	0.43	0.0282 \pm 0.0002	0.0531 \pm 0.0022	20.00 \pm 0.00	125	{125}
	5 000 (10 000)	100	2.38	0.51	0.0283 \pm 0.0002	0.0544 \pm 0.0017	20.00 \pm 0.00	250	{250}
	10 000 (10 000)	100	3.14	0.64	0.0280 \pm 0.0001	0.0520 \pm 0.0014	20.00 \pm 0.00	500	{500}

Table 7: Experimental results relating to the artificial data of Type I

	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	0.60	0.12	0.0178 \pm 0.0006	0.0641 \pm 0.0047	1.00		
	2 500 (10 000)	100	2.29	0.21	0.0168 \pm 0.0004	0.0559 \pm 0.0033	1.00		
	5 000 (10 000)	100	8.37	0.25	0.0165 \pm 0.0003	0.0531 \pm 0.0026	1.00		
	10 000 (10 000)	100	31.61	0.44	0.0163 \pm 0.0003	0.0511 \pm 0.0033	1.00		
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
VP-SVM $r = 0.1$	1 000 (10 000)	100	1.55	0.29	0.0157 \pm 0.0004	0.0427 \pm 0.0040	16.08 \pm 2.50	62	[22, 114]
	2 500 (10 000)	100	1.66	0.36	0.0146 \pm 0.0002	0.0288 \pm 0.0028	15.79 \pm 2.40	158	[51, 292]
	5 000 (10 000)	100	1.88	0.43	0.0142 \pm 0.0001	0.0208 \pm 0.0016	15.95 \pm 2.53	313	[105, 549]
	10 000 (10 000)	100	3.09	0.59	0.0140 \pm 0.0001	0.0163 \pm 0.0013	15.84 \pm 2.50	631	[231, 1043]
VP-SVM $r = 0.25$	1 000 (10 000)	100	0.58	0.17	0.0156 \pm 0.0004	0.0426 \pm 0.0045	6.21 \pm 0.83	161	[60, 267]
	2 500 (10 000)	100	0.77	0.22	0.0147 \pm 0.0002	0.0299 \pm 0.0030	6.29 \pm 0.81	397	[133, 676]
	5 000 (10 000)	100	1.64	0.30	0.0144 \pm 0.0001	0.0247 \pm 0.0023	6.22 \pm 0.82	804	[285, 1299]
	10 000 (10 000)	100	4.99	0.44	0.0141 \pm 0.0001	0.0194 \pm 0.0016	6.38 \pm 0.80	1567	[597, 2535]
VP-SVM $r = 0.5$	1 000 (10 000)	100	0.36	0.12	0.0159 \pm 0.0004	0.0465 \pm 0.0050	3.47 \pm 0.50	288	[122, 482]
	2 500 (10 000)	100	0.81	0.18	0.0151 \pm 0.0002	0.0370 \pm 0.0031	3.52 \pm 0.50	710	[297, 1248]
	5 000 (10 000)	100	2.52	0.28	0.0150 \pm 0.0002	0.0346 \pm 0.0032	3.47 \pm 0.50	1441	[604, 2508]
	10 000 (10 000)	100	9.36	0.43	0.0147 \pm 0.0002	0.0301 \pm 0.0034	3.58 \pm 0.50	2793	[1249, 4962]
VP-SVM $r = 1$	1 000 (10 000)	100	0.33	0.09	0.0169 \pm 0.0005	0.0562 \pm 0.0047	2.00 \pm 0.00	500	[265, 735]
	2 500 (10 000)	100	1.11	0.14	0.0160 \pm 0.0004	0.0477 \pm 0.0041	2.00 \pm 0.00	1250	[613, 1887]
	5 000 (10 000)	100	4.11	0.23	0.0158 \pm 0.0004	0.0452 \pm 0.0044	2.00 \pm 0.00	2500	[1232, 3768]
	10 000 (10 000)	100	16.26	0.42	0.0154 \pm 0.0005	0.0402 \pm 0.0064	2.00 \pm 0.00	5000	[2533, 7467]
VP-SVM $r = 2$	1 000 (10 000)	100	0.37	0.08	0.0178 \pm 0.0006	0.0641 \pm 0.0045	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	2.03	0.12	0.0168 \pm 0.0003	0.0558 \pm 0.0033	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	8.12	0.20	0.0165 \pm 0.0003	0.0531 \pm 0.0026	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	31.12	0.36	0.0163 \pm 0.0003	0.0511 \pm 0.0033	1.00 \pm 0.00	10 000	{10000}
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
RC-SVM $\# \text{ ws} = 1$	1 000 (10 000)	100	0.37	0.07	0.0178 \pm 0.0006	0.0640 \pm 0.0046	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	2.02	0.13	0.0168 \pm 0.0003	0.0558 \pm 0.0033	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	8.14	0.19	0.0165 \pm 0.0003	0.0531 \pm 0.0026	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	31.45	0.36	0.0163 \pm 0.0003	0.0511 \pm 0.0033	1.00 \pm 0.00	10 000	{10000}
RC-SVM $\# \text{ ws} = 2$	1 000 (10 000)	100	0.30	0.11	0.0172 \pm 0.0005	0.0595 \pm 0.0046	2.00 \pm 0.00	500	{500}
	2 500 (10 000)	100	1.10	0.15	0.0169 \pm 0.0004	0.0567 \pm 0.0037	2.00 \pm 0.00	1250	{1250}
	5 000 (10 000)	100	4.06	0.25	0.0165 \pm 0.0003	0.0528 \pm 0.0029	2.00 \pm 0.00	2500	{2500}
	10 000 (10 000)	100	16.33	0.41	0.0163 \pm 0.0002	0.0516 \pm 0.0020	2.00 \pm 0.00	5000	{5000}
RC-SVM $\# \text{ ws} = 3$	1 000 (10 000)	100	0.32	0.14	0.0168 \pm 0.0004	0.0555 \pm 0.0040	3.00 \pm 0.00	333	[333, 334]
	2 500 (10 000)	100	0.88	0.17	0.0169 \pm 0.0004	0.0566 \pm 0.0033	3.00 \pm 0.00	833	[833, 834]
	5 000 (10 000)	100	2.71	0.28	0.0165 \pm 0.0002	0.0534 \pm 0.0022	3.00 \pm 0.00	1667	[1666, 1667]
	10 000 (10 000)	100	10.99	0.45	0.0163 \pm 0.0002	0.0515 \pm 0.0018	3.00 \pm 0.00	3333	[3333, 3334]
RC-SVM $\# \text{ ws} = 4$	1 000 (10 000)	100	0.39	0.17	0.0166 \pm 0.0004	0.0544 \pm 0.0039	4.00 \pm 0.00	250	{250}
	2 500 (10 000)	100	0.81	0.21	0.0166 \pm 0.0003	0.0537 \pm 0.0032	4.00 \pm 0.00	625	{625}
	5 000 (10 000)	100	2.23	0.31	0.0165 \pm 0.0003	0.0532 \pm 0.0027	4.00 \pm 0.00	1250	{1250}
	10 000 (10 000)	100	8.16	0.49	0.0163 \pm 0.0002	0.0516 \pm 0.0020	4.00 \pm 0.00	2500	{2500}
RC-SVM $\# \text{ ws} = 5$	1 000 (10 000)	100	0.46	0.19	0.0167 \pm 0.0004	0.0553 \pm 0.0041	5.00 \pm 0.00	200	{200}
	2 500 (10 000)	100	0.78	0.25	0.0165 \pm 0.0004	0.0532 \pm 0.0035	5.00 \pm 0.00	500	{500}
	5 000 (10 000)	100	1.94	0.32	0.0164 \pm 0.0003	0.0524 \pm 0.0029	5.00 \pm 0.00	1000	{1000}
	10 000 (10 000)	100	6.51	0.51	0.0164 \pm 0.0002	0.0521 \pm 0.0022	5.00 \pm 0.00	2000	{2000}
RC-SVM $\# \text{ ws} = 6$	1 000 (10 000)	100	0.53	0.21	0.0168 \pm 0.0004	0.0556 \pm 0.0031	6.00 \pm 0.00	167	[166, 167]
	2 500 (10 000)	100	0.76	0.28	0.0164 \pm 0.0003	0.0520 \pm 0.0029	6.00 \pm 0.00	417	[416, 417]
	5 000 (10 000)	100	1.77	0.37	0.0164 \pm 0.0003	0.0518 \pm 0.0027	6.00 \pm 0.00	833	[833, 834]
	10 000 (10 000)	100	5.49	0.54	0.0165 \pm 0.0002	0.0530 \pm 0.0019	6.00 \pm 0.00	1667	[1666, 1667]
RC-SVM $\# \text{ ws} = 10$	1 000 (10 000)	100	0.85	0.27	0.0171 \pm 0.0003	0.0585 \pm 0.0027	10.00 \pm 0.00	100	{100}
	2 500 (10 000)	100	0.98	0.40	0.0160 \pm 0.0003	0.0487 \pm 0.0027	10.00 \pm 0.00	250	{250}
	5 000 (10 000)	100	1.56	0.52	0.0160 \pm 0.0002	0.0478 \pm 0.0023	10.00 \pm 0.00	500	{500}
	10 000 (10 000)	100	3.91	0.63	0.0164 \pm 0.0002	0.0524 \pm 0.0017	10.00 \pm 0.00	1000	{1000}
RC-SVM $\# \text{ ws} = 20$	1 000 (10 000)	100	1.83	0.37	0.0196 \pm 0.0007	0.0774 \pm 0.0042	20.00 \pm 0.00	50	{50}
	2 500 (10 000)	100	1.97	0.56	0.0163 \pm 0.0002	0.0512 \pm 0.0020	20.00 \pm 0.00	125	{125}
	5 000 (10 000)	100	2.20	0.79	0.0156 \pm 0.0002	0.0446 \pm 0.0018	20.00 \pm 0.00	250	{250}
	10 000 (10 000)	100	3.19	0.99	0.0160 \pm 0.0002	0.0488 \pm 0.0018	20.00 \pm 0.00	500	{500}

	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	0.36	0.05	0.0588 \pm 0.0009	0.0799 \pm 0.0056	1.00		
	2 500 (10 000)	100	1.82	0.07	0.0572 \pm 0.0005	0.0677 \pm 0.0039	1.00		
	5 000 (10 000)	100	7.47	0.11	0.0580 \pm 0.0003	0.0745 \pm 0.0019	1.00		
	10 000 (10 000)	100	29.56	0.24	0.0559 \pm 0.0003	0.0563 \pm 0.0028	1.00		
VP-SVM $r = 0.1$	1 000 (10 000)	100	1.52	0.29	0.0580 \pm 0.0009	0.0731 \pm 0.0057	15.55 \pm 2.36	64	[26, 109]
	2 500 (10 000)	100	1.61	0.35	0.0553 \pm 0.0004	0.0493 \pm 0.0038	15.99 \pm 2.39	156	[56, 279]
	5 000 (10 000)	100	1.61	0.42	0.0539 \pm 0.0002	0.0347 \pm 0.0024	15.80 \pm 2.45	316	[116, 531]
	10 000 (10 000)	100	3.14	0.58	0.0534 \pm 0.0001	0.0257 \pm 0.0016	15.73 \pm 2.58	636	[236, 1068]
VP-SVM $r = 0.25$	1 000 (10 000)	100	0.53	0.16	0.0563 \pm 0.0007	0.0600 \pm 0.0057	6.18 \pm 0.87	162	[56, 259]
	2 500 (10 000)	100	0.77	0.19	0.0547 \pm 0.0003	0.0429 \pm 0.0034	6.26 \pm 0.81	399	[150, 658]
	5 000 (10 000)	100	1.65	0.26	0.0537 \pm 0.0001	0.0315 \pm 0.0024	6.27 \pm 0.81	797	[321, 1270]
	10 000 (10 000)	100	5.11	0.39	0.0534 \pm 0.0001	0.0244 \pm 0.0021	6.34 \pm 0.79	1577	[597, 2534]
VP-SVM $r = 0.5$	1 000 (10 000)	100	0.34	0.11	0.0564 \pm 0.0007	0.0620 \pm 0.0057	3.48 \pm 0.50	287	[126, 501]
	2 500 (10 000)	100	0.81	0.14	0.0551 \pm 0.0004	0.0478 \pm 0.0047	3.49 \pm 0.50	716	[297, 1244]
	5 000 (10 000)	100	2.44	0.22	0.0543 \pm 0.0004	0.0409 \pm 0.0048	3.50 \pm 0.50	1429	[641, 2517]
	10 000 (10 000)	100	9.40	0.30	0.0538 \pm 0.0004	0.0322 \pm 0.0054	3.54 \pm 0.50	2825	[1209, 4993]
VP-SVM $r = 1$	1 000 (10 000)	100	0.28	0.08	0.0579 \pm 0.0009	0.0738 \pm 0.0062	2.00 \pm 0.00	500	[248, 752]
	2 500 (10 000)	100	1.04	0.11	0.0565 \pm 0.0007	0.0617 \pm 0.0059	2.00 \pm 0.00	1250	[625, 1875]
	5 000 (10 000)	100	3.92	0.15	0.0564 \pm 0.0005	0.0622 \pm 0.0044	2.00 \pm 0.00	2500	[1281, 3719]
	10 000 (10 000)	100	15.93	0.23	0.0555 \pm 0.0004	0.0520 \pm 0.0035	2.00 \pm 0.00	5000	[2404, 7596]
VP-SVM $r = 2$	1 000 (10 000)	100	0.35	0.06	0.0588 \pm 0.0009	0.0799 \pm 0.0057	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	1.84	0.08	0.0572 \pm 0.0005	0.0677 \pm 0.0039	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	7.51	0.12	0.0580 \pm 0.0003	0.0745 \pm 0.0019	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	29.75	0.25	0.0559 \pm 0.0003	0.0563 \pm 0.0028	1.00 \pm 0.00	10000	{10000}
RC-SVM # ws = 1	1 000 (10 000)	100	0.35	0.07	0.0588 \pm 0.0009	0.0801 \pm 0.0058	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	1.85	0.08	0.0572 \pm 0.0005	0.0678 \pm 0.0039	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	7.73	0.12	0.0580 \pm 0.0003	0.0745 \pm 0.0019	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	30.24	0.27	0.0559 \pm 0.0003	0.0563 \pm 0.0028	1.00 \pm 0.00	10000	{10000}
RC-SVM # ws = 2	1 000 (10 000)	100	0.30	0.08	0.0596 \pm 0.0009	0.0855 \pm 0.0052	2.00 \pm 0.00	500	{500}
	2 500 (10 000)	100	1.01	0.11	0.0579 \pm 0.0005	0.0736 \pm 0.0033	2.00 \pm 0.00	1250	{1250}
	5 000 (10 000)	100	3.81	0.15	0.0573 \pm 0.0005	0.0699 \pm 0.0038	2.00 \pm 0.00	2500	{2500}
	10 000 (10 000)	100	15.22	0.24	0.0562 \pm 0.0004	0.0598 \pm 0.0035	2.00 \pm 0.00	5000	{5000}
RC-SVM # ws = 3	1 000 (10 000)	100	0.32	0.12	0.0591 \pm 0.0010	0.0825 \pm 0.0067	3.00 \pm 0.00	333	{333, 334}
	2 500 (10 000)	100	0.82	0.14	0.0585 \pm 0.0006	0.0776 \pm 0.0038	3.00 \pm 0.00	833	{833, 834}
	5 000 (10 000)	100	2.52	0.18	0.0574 \pm 0.0004	0.0710 \pm 0.0031	3.00 \pm 0.00	1667	{1666, 1667}
	10 000 (10 000)	100	10.05	0.24	0.0564 \pm 0.0004	0.0622 \pm 0.0032	3.00 \pm 0.00	3333	{3333, 3334}
RC-SVM # ws = 4	1 000 (10 000)	100	0.37	0.15	0.0589 \pm 0.0011	0.0814 \pm 0.0069	4.00 \pm 0.00	250	{250}
	2 500 (10 000)	100	0.77	0.16	0.0587 \pm 0.0006	0.0793 \pm 0.0038	4.00 \pm 0.00	625	{625}
	5 000 (10 000)	100	2.07	0.20	0.0579 \pm 0.0004	0.0747 \pm 0.0025	4.00 \pm 0.00	1250	{1250}
	10 000 (10 000)	100	7.41	0.27	0.0563 \pm 0.0004	0.0616 \pm 0.0036	4.00 \pm 0.00	2500	{2500}
RC-SVM # ws = 5	1 000 (10 000)	100	0.46	0.17	0.0588 \pm 0.0010	0.0811 \pm 0.0064	5.00 \pm 0.00	200	{200}
	2 500 (10 000)	100	0.73	0.19	0.0586 \pm 0.0007	0.0787 \pm 0.0045	5.00 \pm 0.00	500	{500}
	5 000 (10 000)	100	1.87	0.24	0.0583 \pm 0.0004	0.0775 \pm 0.0027	5.00 \pm 0.00	1000	{1000}
	10 000 (10 000)	100	5.88	0.30	0.0563 \pm 0.0003	0.0617 \pm 0.0029	5.00 \pm 0.00	2000	{2000}
RC-SVM # ws = 6	1 000 (10 000)	100	0.52	0.19	0.0586 \pm 0.0010	0.0796 \pm 0.0070	6.00 \pm 0.00	167	{166, 167}
	2 500 (10 000)	100	0.73	0.21	0.0586 \pm 0.0006	0.0788 \pm 0.0044	6.00 \pm 0.00	417	{416, 417}
	5 000 (10 000)	100	1.67	0.26	0.0586 \pm 0.0005	0.0793 \pm 0.0033	6.00 \pm 0.00	833	{833, 834}
	10 000 (10 000)	100	5.01	0.33	0.0565 \pm 0.0003	0.0637 \pm 0.0026	6.00 \pm 0.00	1667	{1666, 1667}
RC-SVM # ws = 10	1 000 (10 000)	100	0.89	0.25	0.0593 \pm 0.0011	0.0845 \pm 0.0069	10.00 \pm 0.00	100	{100}
	2 500 (10 000)	100	0.99	0.33	0.0579 \pm 0.0006	0.0744 \pm 0.0043	10.00 \pm 0.00	250	{250}
	5 000 (10 000)	100	1.51	0.36	0.0590 \pm 0.0005	0.0824 \pm 0.0033	10.00 \pm 0.00	500	{500}
	10 000 (10 000)	100	3.59	0.44	0.0574 \pm 0.0003	0.0705 \pm 0.0024	10.00 \pm 0.00	1000	{1000}
RC-SVM # ws = 20	1 000 (10 000)	100	1.90	0.38	0.0671 \pm 0.0019	0.1238 \pm 0.0082	20.00 \pm 0.00	50	{50}
	2 500 (10 000)	100	2.01	0.53	0.0578 \pm 0.0007	0.0742 \pm 0.0052	20.00 \pm 0.00	125	{125}
	5 000 (10 000)	100	2.19	0.64	0.0583 \pm 0.0005	0.0779 \pm 0.0031	20.00 \pm 0.00	250	{250}
	10 000 (10 000)	100	3.20	0.67	0.0582 \pm 0.0003	0.0770 \pm 0.0021	20.00 \pm 0.00	500	{500}

Table 9: Experimental results relating to the artificial data of Type III

	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	0.61	0.07	0.0155 \pm 0.0002	0.0848 \pm 0.0012	1.00		
	2 500 (10 000)	100	3.28	0.16	0.0151 \pm 0.0003	0.0821 \pm 0.0016	1.00		
	5 000 (10 000)	100	11.89	0.63	0.0132 \pm 0.0007	0.0694 \pm 0.0049	1.00		
	10 000 (10 000)	100	45.94	0.58	0.0137 \pm 0.0007	0.0730 \pm 0.0046	1.00		
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
VP-SVM $r = 0.25$	1 000 (10 000)	100	4.27	0.62	0.0172 \pm 0.0005	0.0946 \pm 0.0024	44.23 \pm 2.19	23	[4, 54]
	2 500 (10 000)	100	4.78	0.85	0.0142 \pm 0.0002	0.0775 \pm 0.0015	46.09 \pm 2.58	54	[13, 108]
	5 000 (10 000)	100	5.03	1.16	0.0129 \pm 0.0002	0.0682 \pm 0.0010	46.32 \pm 2.61	108	[19, 244]
	10 000 (10 000)	100	5.72	1.76	0.0119 \pm 0.0001	0.0607 \pm 0.0007	47.52 \pm 2.72	210	[43, 416]
VP-SVM $r = 0.5$	1 000 (10 000)	100	1.45	0.30	0.0158 \pm 0.0003	0.0867 \pm 0.0019	14.50 \pm 1.25	69	[18, 155]
	2 500 (10 000)	100	1.62	0.49	0.0139 \pm 0.0002	0.0752 \pm 0.0013	15.22 \pm 1.47	164	[40, 323]
	5 000 (10 000)	100	2.17	0.78	0.0129 \pm 0.0002	0.0683 \pm 0.0011	15.00 \pm 1.45	333	[78, 635]
	10 000 (10 000)	100	4.52	1.24	0.0120 \pm 0.0001	0.0611 \pm 0.0009	15.72 \pm 1.62	636	[144, 1261]
VP-SVM $r = 1$	1 000 (10 000)	100	0.49	0.18	0.0152 \pm 0.0003	0.0833 \pm 0.0017	5.16 \pm 0.53	194	[69, 478]
	2 500 (10 000)	100	0.99	0.32	0.0140 \pm 0.0003	0.0753 \pm 0.0019	5.40 \pm 0.51	463	[156, 1260]
	5 000 (10 000)	100	2.77	0.58	0.0132 \pm 0.0003	0.0697 \pm 0.0020	5.48 \pm 0.54	912	[264, 2437]
	10 000 (10 000)	100	10.03	0.80	0.0124 \pm 0.0002	0.0638 \pm 0.0013	5.56 \pm 0.64	1799	[529, 4868]
VP-SVM $r = 2$	1 000 (10 000)	100	0.43	0.10	0.0155 \pm 0.0002	0.0847 \pm 0.0014	1.72 \pm 0.45	581	[283, 1000]
	2 500 (10 000)	100	2.29	0.15	0.0149 \pm 0.0004	0.0808 \pm 0.0024	1.55 \pm 0.50	1613	[696, 2500]
	5 000 (10 000)	100	7.96	0.42	0.0138 \pm 0.0007	0.0739 \pm 0.0049	1.67 \pm 0.47	2994	[1345, 5000]
	10 000 (10 000)	100	31.92	0.82	0.0126 \pm 0.0009	0.0656 \pm 0.0065	1.64 \pm 0.48	6098	[2702, 10000]
VP-SVM $r = 3$	1 000 (10 000)	100	0.54	0.07	0.0155 \pm 0.0002	0.0847 \pm 0.0012	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	3.09	0.12	0.0151 \pm 0.0003	0.0822 \pm 0.0016	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	11.64	0.53	0.0132 \pm 0.0007	0.0694 \pm 0.0049	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	45.40	0.54	0.0137 \pm 0.0007	0.0730 \pm 0.0046	1.00 \pm 0.00	10000	{10000}
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
RC-SVM # ws = 1	1 000 (10 000)	100	0.53	0.06	0.0155 \pm 0.0002	0.0847 \pm 0.0012	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	3.08	0.12	0.0151 \pm 0.0003	0.0822 \pm 0.0016	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	11.67	0.52	0.0132 \pm 0.0007	0.0694 \pm 0.0049	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	46.22	0.55	0.0137 \pm 0.0007	0.0730 \pm 0.0046	1.00 \pm 0.00	10000	{10000}
RC-SVM # ws = 2	1 000 (10 000)	100	0.37	0.14	0.0153 \pm 0.0002	0.0837 \pm 0.0015	2.00 \pm 0.00	500	{500}
	2 500 (10 000)	100	1.61	0.13	0.0152 \pm 0.0002	0.0830 \pm 0.0012	2.00 \pm 0.00	1250	{1250}
	5 000 (10 000)	100	6.12	0.22	0.0151 \pm 0.0002	0.0825 \pm 0.0014	2.00 \pm 0.00	2500	{2500}
	10 000 (10 000)	100	23.51	0.97	0.0130 \pm 0.0006	0.0680 \pm 0.0046	2.00 \pm 0.00	5000	{5000}
RC-SVM # ws = 5	1 000 (10 000)	100	0.46	0.22	0.0152 \pm 0.0002	0.0830 \pm 0.0014	5.00 \pm 0.00	200	{200}
	2 500 (10 000)	100	0.96	0.30	0.0150 \pm 0.0002	0.0815 \pm 0.0009	5.00 \pm 0.00	500	{500}
	5 000 (10 000)	100	2.70	0.28	0.0153 \pm 0.0001	0.0833 \pm 0.0008	5.00 \pm 0.00	1000	{1000}
	10 000 (10 000)	100	9.95	0.36	0.0152 \pm 0.0001	0.0826 \pm 0.0007	5.00 \pm 0.00	2000	{2000}
RC-SVM # ws = 10	1 000 (10 000)	100	0.88	0.29	0.0154 \pm 0.0002	0.0844 \pm 0.0011	10.00 \pm 0.00	100	{100}
	2 500 (10 000)	100	1.02	0.46	0.0148 \pm 0.0001	0.0808 \pm 0.0007	10.00 \pm 0.00	250	{250}
	5 000 (10 000)	100	1.96	0.58	0.0150 \pm 0.0001	0.0818 \pm 0.0008	10.00 \pm 0.00	500	{500}
	10 000 (10 000)	100	5.67	0.53	0.0152 \pm 0.0001	0.0827 \pm 0.0005	10.00 \pm 0.00	1000	{1000}
RC-SVM # ws = 15	1 000 (10 000)	100	1.50	0.34	0.0157 \pm 0.0002	0.0864 \pm 0.0014	15.00 \pm 0.00	67	[66, 67]
	2 500 (10 000)	100	1.62	0.53	0.0149 \pm 0.0001	0.0814 \pm 0.0007	15.00 \pm 0.00	167	[166, 167]
	5 000 (10 000)	100	2.10	0.84	0.0149 \pm 0.0001	0.0812 \pm 0.0006	15.00 \pm 0.00	333	[333, 334]
	10 000 (10 000)	100	4.80	0.87	0.0150 \pm 0.0001	0.0819 \pm 0.0005	15.00 \pm 0.00	667	[666, 667]
RC-SVM # ws = 20	1 000 (10 000)	100	1.88	0.39	0.0162 \pm 0.0003	0.0891 \pm 0.0017	20.00 \pm 0.00	50	{50}
	2 500 (10 000)	100	1.99	0.59	0.0151 \pm 0.0001	0.0822 \pm 0.0006	20.00 \pm 0.00	125	{125}
	5 000 (10 000)	100	2.27	0.89	0.0149 \pm 0.0001	0.0814 \pm 0.0006	20.00 \pm 0.00	250	{250}
	10 000 (10 000)	100	4.24	1.14	0.0149 \pm 0.0001	0.0812 \pm 0.0005	20.00 \pm 0.00	500	{500}
RC-SVM # ws = 40	1 000 (10 000)	100	3.96	0.59	0.0194 \pm 0.0007	0.1053 \pm 0.0033	40.00 \pm 0.00	25	{25}
	2 500 (10 000)	100	4.06	0.82	0.0157 \pm 0.0002	0.0859 \pm 0.0010	40.00 \pm 0.00	62	[62, 63]
	5 000 (10 000)	100	4.29	1.17	0.0152 \pm 0.0001	0.0829 \pm 0.0005	40.00 \pm 0.00	125	{125}
	10 000 (10 000)	100	4.93	1.75	0.0148 \pm 0.0001	0.0808 \pm 0.0004	40.00 \pm 0.00	250	{250}
RC-SVM # ws = 50	1 000 (10 000)	100	5.06	0.69	0.0222 \pm 0.0011	0.1181 \pm 0.0045	50.00 \pm 0.00	20	{20}
	2 500 (10 000)	100	4.99	0.93	0.0161 \pm 0.0002	0.0883 \pm 0.0011	50.00 \pm 0.00	50	{50}
	5 000 (10 000)	100	5.30	1.29	0.0153 \pm 0.0001	0.0839 \pm 0.0005	50.00 \pm 0.00	100	{100}
	10 000 (10 000)	100	6.43	1.94	0.0149 \pm 0.0001	0.0810 \pm 0.0004	50.00 \pm 0.00	200	{200}

	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
LS-SVM	1 000 (10 000)	100	0.53	0.05	0.0649 \pm 0.0004	0.0370 \pm 0.0055	1.00		
	2 500 (10 000)	100	3.26	0.09	0.0647 \pm 0.0002	0.0330 \pm 0.0032	1.00		
	5 000 (10 000)	100	12.51	0.15	0.0640 \pm 0.0001	0.0240 \pm 0.0015	1.00		
	10 000 (10 000)	100	49.29	0.28	0.0640 \pm 0.0001	0.0223 \pm 0.0015	1.00		
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
VP-SVM $r = 0.25$	1 000 (10 000)	100	4.85	0.66	0.0797 \pm 0.0021	0.1258 \pm 0.0076	44.56 \pm 2.07	22	[6 , 54]
	2 500 (10 000)	100	4.76	0.77	0.0719 \pm 0.0009	0.0882 \pm 0.0043	46.18 \pm 2.55	54	[14 , 116]
	5 000 (10 000)	100	4.88	0.95	0.0682 \pm 0.0005	0.0663 \pm 0.0033	47.45 \pm 2.61	105	[20 , 214]
	10 000 (10 000)	100	5.77	1.14	0.0657 \pm 0.0003	0.0460 \pm 0.0027	48.05 \pm 2.78	208	[38 , 407]
VP-SVM $r = 0.5$	1 000 (10 000)	100	1.65	0.30	0.0703 \pm 0.0011	0.0820 \pm 0.0068	14.98 \pm 1.44	67	[16 , 129]
	2 500 (10 000)	100	1.64	0.37	0.0670 \pm 0.0006	0.0565 \pm 0.0045	15.17 \pm 1.52	165	[40 , 334]
	5 000 (10 000)	100	2.26	0.45	0.0652 \pm 0.0003	0.0404 \pm 0.0035	15.69 \pm 1.47	319	[74 , 644]
	10 000 (10 000)	100	4.81	0.62	0.0645 \pm 0.0002	0.0307 \pm 0.0023	15.83 \pm 1.49	632	[167 , 1343]
VP-SVM $r = 1$	1 000 (10 000)	100	0.56	0.16	0.0670 \pm 0.0008	0.0578 \pm 0.0064	5.33 \pm 0.47	188	[68 , 464]
	2 500 (10 000)	100	1.05	0.19	0.0656 \pm 0.0004	0.0431 \pm 0.0042	5.33 \pm 0.47	469	[171 , 1162]
	5 000 (10 000)	100	3.01	0.26	0.0643 \pm 0.0002	0.0290 \pm 0.0025	5.42 \pm 0.52	923	[297 , 2445]
	10 000 (10 000)	100	10.81	0.41	0.0641 \pm 0.0001	0.0243 \pm 0.0022	5.54 \pm 0.61	1805	[555 , 4840]
VP-SVM $r = 2$	1 000 (10 000)	100	0.49	0.13	0.0653 \pm 0.0006	0.0427 \pm 0.0061	1.65 \pm 0.48	606	[266 , 1000]
	2 500 (10 000)	100	2.21	0.11	0.0649 \pm 0.0003	0.0357 \pm 0.0042	1.70 \pm 0.48	1471	[603 , 2500]
	5 000 (10 000)	100	9.32	0.18	0.0641 \pm 0.0001	0.0248 \pm 0.0018	1.54 \pm 0.50	3247	[1348 , 5000]
	10 000 (10 000)	100	35.16	0.31	0.0640 \pm 0.0001	0.0227 \pm 0.0018	1.62 \pm 0.49	6173	[2608 , 10000]
VP-SVM $r = 3$	1 000 (10 000)	100	0.54	0.07	0.0649 \pm 0.0004	0.0371 \pm 0.0055	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	3.25	0.09	0.0647 \pm 0.0002	0.0331 \pm 0.0032	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	12.54	0.16	0.0640 \pm 0.0001	0.0240 \pm 0.0015	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	49.35	0.29	0.0640 \pm 0.0001	0.0223 \pm 0.0015	1.00 \pm 0.00	10 000	{10000}
	data set sizes	runs	train time	test time	test error	L_2 -error	# of ws	ws size: median	ws size: range
RC-SVM # ws = 1	1 000 (10 000)	100	0.53	0.07	0.0649 \pm 0.0004	0.0373 \pm 0.0056	1.00 \pm 0.00	1000	{1000}
	2 500 (10 000)	100	3.27	0.09	0.0647 \pm 0.0002	0.0331 \pm 0.0033	1.00 \pm 0.00	2500	{2500}
	5 000 (10 000)	100	12.55	0.16	0.0640 \pm 0.0001	0.0240 \pm 0.0015	1.00 \pm 0.00	5000	{5000}
	10 000 (10 000)	100	50.04	0.29	0.0640 \pm 0.0001	0.0223 \pm 0.0015	1.00 \pm 0.00	10 000	{10000}
RC-SVM # ws = 2	1 000 (10 000)	100	0.38	0.09	0.0652 \pm 0.0005	0.0409 \pm 0.0056	2.00 \pm 0.00	500	{500}
	2 500 (10 000)	100	1.71	0.12	0.0649 \pm 0.0003	0.0352 \pm 0.0037	2.00 \pm 0.00	1250	{1250}
	5 000 (10 000)	100	6.39	0.18	0.0640 \pm 0.0001	0.0242 \pm 0.0019	2.00 \pm 0.00	2500	{2500}
	10 000 (10 000)	100	25.47	0.30	0.0640 \pm 0.0001	0.0221 \pm 0.0014	2.00 \pm 0.00	5000	{5000}
RC-SVM # ws = 5	1 000 (10 000)	100	0.48	0.17	0.0655 \pm 0.0005	0.0440 \pm 0.0056	5.00 \pm 0.00	200	{200}
	2 500 (10 000)	100	1.02	0.19	0.0654 \pm 0.0004	0.0426 \pm 0.0044	5.00 \pm 0.00	500	{500}
	5 000 (10 000)	100	2.82	0.24	0.0643 \pm 0.0002	0.0295 \pm 0.0026	5.00 \pm 0.00	1000	{1000}
	10 000 (10 000)	100	10.63	0.37	0.0641 \pm 0.0001	0.0248 \pm 0.0017	5.00 \pm 0.00	2000	{2000}
RC-SVM # ws = 10	1 000 (10 000)	100	0.85	0.25	0.0657 \pm 0.0005	0.0471 \pm 0.0054	10.00 \pm 0.00	100	{100}
	2 500 (10 000)	100	1.05	0.33	0.0656 \pm 0.0004	0.0448 \pm 0.0041	10.00 \pm 0.00	250	{250}
	5 000 (10 000)	100	2.04	0.35	0.0646 \pm 0.0002	0.0335 \pm 0.0032	10.00 \pm 0.00	500	{500}
	10 000 (10 000)	100	6.34	0.51	0.0644 \pm 0.0001	0.0299 \pm 0.0019	10.00 \pm 0.00	1000	{1000}
RC-SVM # ws = 15	1 000 (10 000)	100	1.47	0.32	0.0664 \pm 0.0006	0.0537 \pm 0.0059	15.00 \pm 0.00	67	[66 , 67]
	2 500 (10 000)	100	1.63	0.43	0.0656 \pm 0.0004	0.0449 \pm 0.0040	15.00 \pm 0.00	167	[166 , 167]
	5 000 (10 000)	100	2.18	0.49	0.0647 \pm 0.0002	0.0350 \pm 0.0030	15.00 \pm 0.00	333	[333 , 334]
	10 000 (10 000)	100	5.23	0.58	0.0646 \pm 0.0001	0.0335 \pm 0.0022	15.00 \pm 0.00	667	[666 , 667]
RC-SVM # ws = 20	1 000 (10 000)	100	1.92	0.38	0.0671 \pm 0.0007	0.0601 \pm 0.0054	20.00 \pm 0.00	50	{50}
	2 500 (10 000)	100	2.00	0.52	0.0659 \pm 0.0004	0.0473 \pm 0.0037	20.00 \pm 0.00	125	{125}
	5 000 (10 000)	100	2.32	0.61	0.0647 \pm 0.0002	0.0356 \pm 0.0029	20.00 \pm 0.00	250	{250}
	10 000 (10 000)	100	4.75	0.69	0.0647 \pm 0.0002	0.0353 \pm 0.0021	20.00 \pm 0.00	500	{500}
RC-SVM # ws = 40	1 000 (10 000)	100	3.89	0.60	0.0707 \pm 0.0011	0.0850 \pm 0.0065	40.00 \pm 0.00	25	{25}
	2 500 (10 000)	100	4.05	0.81	0.0669 \pm 0.0005	0.0577 \pm 0.0039	40.00 \pm 0.00	62	[62 , 63]
	5 000 (10 000)	100	4.47	1.00	0.0650 \pm 0.0002	0.0388 \pm 0.0030	40.00 \pm 0.00	125	{125}
	10 000 (10 000)	100	5.66	1.24	0.0649 \pm 0.0002	0.0379 \pm 0.0022	40.00 \pm 0.00	250	{250}
RC-SVM # ws = 50	1 000 (10 000)	100	4.70	0.69	0.0726 \pm 0.0010	0.0960 \pm 0.0051	50.00 \pm 0.00	20	{20}
	2 500 (10 000)	100	5.06	0.90	0.0675 \pm 0.0005	0.0627 \pm 0.0038	50.00 \pm 0.00	50	{50}
	5 000 (10 000)	100	5.33	1.15	0.0651 \pm 0.0002	0.0410 \pm 0.0027	50.00 \pm 0.00	100	{100}
	10 000 (10 000)	100	6.57	1.49	0.0649 \pm 0.0002	0.0383 \pm 0.0020	50.00 \pm 0.00	200	{200}

Table 11: Experimental results relating to the artificial data of Type V